AD-753 416

COMPUTATIONAL ALGORITHMS FOR UNCONSTRAINED OPTIMIZATION

Bruce T. Kujawski

Air Force Flight Dynamic Laboratory Wright-Patterson Air Force Base, Ohio

October 1972

DISTRIBUTED BY:



National Technical Information Service
U. S. DEPARTMENT OF COMMERCE

5285 Port Royal Road, Springfield Va. 22151

COMPUTATIONAL ALGORITHMS FOR UNCONSTRAINED OPTIMIZATION

BRUCE T. KUJAWSKI, MAJOR, USAF

TECHNICAL REPORT AFFDL-TR-72-77

OCTOBER 1972



Approved for public release; distribution unlimited.

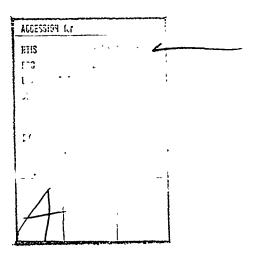
Reproduced by
NATIONAL TECHNICAL
INFORMATION SERVICE
U.S. Department of Commerce

AIR FORCE FLIGHT DYNAMICS LABORATORY
AIR FORCE SYSTEMS COMMAND
WRIGHT-PATTERSON AIR FORCE BASE, OHIO



NOTICE

When Government drawings, specifications, or other data are used for any purpose other than in connection with a definitely related Government procurement operation, the United States Government thereby incurs no responsibility nor any obligation whatsoever; and the fact that the government may have formulated, furnished, or in any way supplied the said drawings, specifications, or other data, is not to be regarded by implication or otherwise as in any manner licensing the holder or any other person or corporation, or conveying any rights or permission to manufacture, use, or sell any patented invention that may in any way be related thereto.



Copies of this report should not be returned unless return is required by security considerations, contractual obligations, or notice on a specific document.

AIR FORCE,56780/7 December 1972 - 300

Security Classification								
DOCUMENT CONTROL DATA - R & D (Security classification of title, body of abatract and indexing annotation must be entered when the overall report is classified)								
1 ORIGINATING ACTIVITY (Corporate author)	mnotation must be el		CURITY CLASSIFICATION					
, , , , ,								
Air Force Flight Dynamics Laboratory Wright-Patterson AFB, Ohio 45433		Unclassified						
		20. GROUP						
3 REPORT TITLE								
COMPUTATIONAL ALGORITAMS FOR UNCONSTRAINED OPTIMIZATION								
4. DESCRIPTIVE NOTES (Type of report and inclusive dates)								
5 AUTHORIS) (First name, middle initial, last name)								
Bruce T. Kujawski, Major, USAF								
6 REPORT DATE	74. TOTAL NO OF	PAGES	7b. NO. OF REFS					
October 1972			15					
NO. CONTRACT OR GRANT NO.	24. ORIGINATOR'S	REPORT NUMBER(S)						
	}							
B. PROJECT NO 8219	AFFDL-TR-72-77							
- Task No. 821911	9b. OTHER REPORT NO(S) (Any other numbers that may be easigned this rep' it)							
d.								
10 DISTRIBUTION STATEMENT								
Approved for public release; distribution unlimited.								
11. SUPPLEMENTARY NOTES	12 SPONSORING MILITARY ACTIVITY							
Air Force Institute of Technology	Air Force Flight Dynamics Laboratory							
DS/EE/72-1	Air Force Systems Command							
DO DU C=1	Wright-Patterson AFB, Ohio 45433							
(1) ARSTOLCT								

A generalized descent algorithm theory is developed for unconstrained minimization problems. Here a descent algorithm is defined as a computational procedure where at each iteration a descent direction is determined and a single-dimensional search is made for the minimum in the descert direction. The theory is shown to be a generalization of the three most common descent algorithms; gradient, conjugate gradient, and Fletcher-Powell.

Execution of the single-dimensional search can be computationally time consuming. Two additional algorithms are presented which reduce or eliminate single-dimensional search time. The first is a modification of Davidon's Variance Algorithm and requires a minimal single-dimensional search. The second is a direct method for minimizing a special class of quadratic functions.

DD FORM . 1473

UNCLASSIFIED

10

Security Classification

UNCLASSIFIED

UNCLASSIFIED Security Classification						
14. KEY WORDS	L!N!		LINK B LINK C			
	ROLE	WT.	ROLE	WT	RCLE	WT
Descent Algorithms						
Function Minimizations						
Rank-One						
Conjugate Gradient						
Fletcher-Powell Method	 					
Davidon Variance						
Algorithms						
						1
				İ		
			<u> </u>			

COMPUTATIONAL ALGORITHMS FOR UNCONSTRAINED OPTIMIZATION

BRUCE T. KUJAWSKI, MAJOR, USAF

Approved for public release; distribution unlimited.

THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY OF THE PROPERTY O

FOPEWORD

This report was originally prepared as a dissertation in partial fulfillment of the requirement for the degree of Doctor of Philosophy in Aerospace Ergineering from the Air Force Institute of Technology, Wright-Patterson AFB, Ohio by Captain Bruce T. Kujawski. Capt. Kujawski's advisor was Lt. Colonel Russell Hannen.

The research was accomplished while the author was assigned to the Air Force Flight Dynamics Laboratory, Wright-Patterson AFB, Chic. The research was conducted under Project 8219, Task 821911, "Flight Control Optimization Techniques".

The author gratefully acknowledges the assistance, suggestions, and stimulations provided by many people in the Laboratory, particularly Mr. Ronald O. Anderson, Capt. Don C. Eckholdt, and Mr. Jerry E. Jenking.

Special acknowledgement is given to Dr. Donald O. Norris, formerly of AFIT, and presently at the Ohio University, for his patience, suggestions, and encouragement.

This technical report has been reviewed and is approved.

Chief, Control Criteria Branch

Flight Cortrol Division

AF Flight Dynamics Laboratory

Strong to the strong of the property of the strong of the

ABSTRACT

A generalized descent algorithm theory is developed for unconstrained minimization problems. Here a descent algorithm is defined as a computational procedure where at each iteration a descent direction is determined and a single-dimensional search is made for the minimum in the descent direction. The theory is shown to be a generalization of the three most common descent algorithm; gradient, conjugate gradient, and Fletcher-Powell.

Execution of the single-dimensional search can be computationally time consuming. Two additional algorithms are presented which reduce or eliminate single-dimensional search time. The first is a modification of Davidon's Variance Algorithm and requires a minimal single-dimensional search. The second is a direct method for minimizing a special class of quadratic functions of the form $1/2||x||^2 + 1/2 k(a - m'x)^2$.

TABLE OF CONTENTS

SECTIO	ON		PAGE
I	INT	RODUCTION	1
	1.	Background	1
	2.	Outline and Preview of Sections	2
II	DES	CENT ALGORITHMS	6
	1.	Notation and Definitions	6
	2.	Rasic Theorem on Generalized Descent Algorithms	7
	3.	Application of the Basic Theorem to Common Descent Algorithms	12
III	A F	RANK-ONE METHOD OF FUNCTION MUNIMIZATION	21
	1.	Basic Rank-One Method	21
	2.	Structure of the Rank-One Minimization Algorithms	23
	3.	Complete Rank-One Algorithm for Function Minimization	29
IV	DIF	RECT APPLICATION OF RANK-ONE	32
	1.	Application to an Aircraft Wing-Root Bending Problem	34
	2.	Derivation of the Direct Rank-One Method	88
	CON	ICLUSIONS	47
REFER	ences	3	51

Preceding page blank

and the property of the state o

SYMBOLS

- belongs to, is a member of; x(X: x is a member of the set X
- + approaches, converges to
- + monotonically decreasing; $f(x_n)+L$: $f(x_n)$ monotonically decreases to the value L
- [y,z] value of the continuous linear functional y operating on the vector z. In a Hilbert space this is the inner product of the vector y and z. Also, a closed interval on the real line, alternate to parenthesis, and references; context will make usage clear.
- ||x|| The norm of the vector x
- The union of sets
- C is a subset of (or equal to)
- matrix or vector transpose except in Section II where f' is the derivative of f.
- RN _ iclidean N-space
- D² second differential operator
- $\{M\}_{i,j}$ element of the $i\frac{th}{t}$ row and $j\frac{th}{t}$ column of the matrix M

SECTION I

INTRODUCTION

1. BACKGROUND

In virtually all fields of the physical sciences, and particularly in engineering, the digital computer is the principal tool used in the solution of complex problems. The speed and flexibility of the computer has in many cases changed the nature of the problems that can be solved, i.e., the solution must not only meet specific constraints but must also be the best or optimal in some specified sense. There are three interrelated tasks in the formulation of such an optimization problem.

First, the physical system or process must be described mathematically or modeled in terms appropriate for computation. Second, the measure of goodness, generally referred to as the cost function, penalty function or payoff function must be defined to adequately describe how one solution compares to another. Finally, computational methods must be applied to find a solution which satisfies the mathematical model and cost function in such a way so as to extract the best or optimum solution.

Generally such a problem can be cast into a constrained optimization problem such as: Find the solution x which minimizes the cost function f(x) while satisfying specific constraints described by g(x)=0. Often the problem can be simplified, conceptually, by adjoining the constraints to the cost function through the use of Lagrange multipliers. Thus the constrained optimization problem is converted

to the following unconstrained problem: Find the solution (x,λ) which extremizes the cost function $F(x,\lambda)=f(x)+\lambda g(x)$. The existence of such Lagrange multipliers is a subject in itself.

Another method of solving the constrained optimization problem is to restrict the problem to a subspace, an approximation to the constraints for example, and considering a related unconstrained problem as an intermediate step in obtaining the solution to the constrained problem. Since the simplified unconstrained problem may have to be solved many times in order to obtain the solution of constrained problem, an efficient method of solving the unconstrained problem is essential. Finally, the solution of the unconstrained problem is often of interest in itself.

The subject of this thesis is the computational methods which may be used to arrive at a minimizing solution to the unconstrained minimization problem. It is tacitly assumed that any constraints are accounted for through the use of Lagrange multipliers or other valid techniques, such as penalty factors.

OUTLINE AND PREVIEW OF SECTIONS

For functions which have a continuous first derivative the most common methods used to minimize the function, i.e., obtain the solution to the unconstrained minimization problem, are the gradient, conjugate gradient, and Fletcher-Powell algorithms. These algorithms are reviewed briefly to illustrate certain common elements. Here it is assumed the function to be minimized is f which is defined for each x in some space X. Further, assume the gradient, g, of f at x also exists: $g(x)=grad \ f(x)$. For each algorithm only the initializations required and the recursive equations are given. Convergence criteria or tests

for convergence although important in computational applications are omitted here in order to emphasize those properties which these algorithms have in common.

Gradient Algorithm:

Initially: choose an arbitrary x_0

Iteratively: set $s_n = -g(x_n)$

choose $\alpha = \alpha_n$ to minimize $f(x_n + \alpha s_n)$

set $x_{n+1}=x_n + \alpha_n s_n$.

Conjugate Gradient Algorithm: (Reference 5)

Initially: choose an arbitrary x₀

 $set s_o = -g(x_o)$

Iteratively: choose $\alpha = \alpha_n$ to minimize $f(x_n + \alpha s_n)$

set $x_{n+1} = x_n + \alpha_n s_n$

 $\beta_n = \frac{||g(x_{n+1})||^2}{||g(x_n)||^2}$

 $s_{n+1} = -g(x_{n+1}) + \beta_n s_n.$

Pletcher-Powell Algorithm: (Reference 4)

Initially: choose an arbitrary x_o

set $H_n = I$

Iteratively: $s_n = -H_ng(x_n)$

choose $a = \alpha_n$ to minimize $f(x_n + \alpha s_n)$

set $x_{n+1} = x_n + \alpha_n s_n$

 $\sigma_n = x_{n+1} - x_n$

 $y_n = g(x_{n+1}) - g(x_n)$

 $\mathbf{H}_{n+1} = \mathbf{H}_{n} - \frac{\mathbf{H}_{n} \mathbf{y}_{n} \mathbf{y}_{n}^{\dagger} \mathbf{H}_{n}}{\mathbf{y}_{n}^{\dagger} \mathbf{H}_{n} \mathbf{y}_{n}} + \frac{\sigma_{n} \sigma_{n}^{\dagger}}{\sigma_{n}^{\dagger} \mathbf{y}_{n}}$

where the prime (') denotes transpose.

AT /TL-TR-72-77

Each of these algorithms generates a search direction, s_n , for which the function, initially at least, tends to decrease, i.e., for which $g'(x_n)s_n$ <0. A single-dimensional search is then conducted to obtain the minimum of f in the direction s_n from the current point x_n . The location of the minimum of the single-dimensional search is chosen as the next iteration point, x_{n+1} . The differences between the algorithms are in the method used to generate the search directions s_n . These algorithms and others which generate a descent direction and incorporate a single-dimensional search will be collectively classed as descent algorithms.

In Section II the proof of a theorem which is a generalization of descent algorithms is presented. Specific applications to the gradient, Fletcher-Powell, and conjugate gradient algorithm are given at the end of the chapter.

Next, consider the problem of minimizing the quadratic function $f(x) = f_0 + a^*x + 1/2x^*Gx$. The gradient of f at x is given by g(x) = a + Gx. Let $h = -G^{-1}g(x)$, assuming G^{-1} exists, then g(x+h) = a + G(x+h) = a + Gx + Gh = g(x) - g(x)

that is, $x^* = x^+h$ satisfies $g(x^*) = 0$, the necessary condition for f to have a minimum at x^* . Note that G is the second derivative of f. Clearly for a quadratic function, knowledge of G or the second derivative, or better G^{-1} , greatly simplifies the problem of finding the minimum of the function. Since many functions can be approximated by a

quadratic in some neightborhood of a (local) minimum, information about the second derivative, or its inverse, should enhance the ability to arrive at a solution to the general unconstrained minimization problem.

The Meston-Raphson algorithm is a method of function minimization which utilizes the inverse of the second derivative. In addition to the computational difficulties of obtaining the second derivative, this method requires an initial estimate sufficiently close to the final solution before convergence is guaranteed. Because of these difficulties, several algorithms termed quasi-Newton methods by Powell (Reference 12 have been constructed which iteratively estimate the inverse of the second derivative. The best known of these is the method of Fletcher-Powell. Another more recent method of this type is Davidon's Variance Algorithm, not to be confused with Davidon's Variance Metric Algorithm which was the predecessor to Fletcher-Powell's method. A new derivation of Dav. don's Variance Algorithm is presented in Section III.

Davidon's algorithm suffers some difficulties and Section III concludes with a modified version of the algorithm which, although somewhat more complex, circumvents one major difficulty.

Whenever a method is available for obtaining the inverse of the second derivative, particularly for a quadratic function, the minimizing solution can be obtained directly. In Section IV a method for the direct solution of a special class of quadratic minimization problems is presented. The procedure is based on the Rank-One method of matrix inversion. The algorithm contains a necessary and sufficient test for the existence of an extremum and a sufficient test that the extremum be a minimum. The special class of problems to which this method applies are generalizations of the following form: $f(x) = 1/2 |x||^2 + k(a-m'x)^2$.

AFFD1-TB-72-77

SECTION II

DESCENT ALGORITHMS

In this section the proof of a basic theorem on descent algorithms is presented followed by applications to several familiar descent algorithms.

1. NOVATION AND DEFINITIONS

Suppose X denotes a real normed linear space and f a real-valued function defined on X. For an arbitrary point x_0 of X, denote by S the "level set" of f at x_0 , i.e., $S = \{x : f(x) \leq f(x_0)\}$. The Frechet derivative of f at x will be denoted f'(x) and if $x^{\pm} \in X^{\pm}$, the topological dual of X, the value of x^{\pm} at x will be written $[x^{\pm}, x]$.

Let \$\phi\$ denote a bounded map from S to X satisfying:

- (i) $[f'(x), \phi(x)] \ge 0$ for all $x \in S$, and
- (ii) given an $\varepsilon>0$ there exists $\delta>0$ such that $[f'(x), \phi(x)]<\delta$ implies $||f'(x)||<\varepsilon$.

Observe for later reference that condition (ii) implies that $f'(x) = 0 \text{ whenever } [f'(x), \phi(x)] = 0 \text{ for if there exists an } x_1 \text{ such that } [f'(x_1), \phi(x_1)] = 0 \text{ but } f'(x_1) \neq 0 \text{ set } \varepsilon = 1/2 ||f'(x_1)|| > 0 \text{ then } [f'(x_1), \phi(x_1)] = 0 < \delta \text{ for all } \delta \text{ while } ||f'(x_1)|| = 2\varepsilon > \varepsilon \text{ contrary to (ii)}.$ Condition (ii) also implies $[f'(x), \phi(x)]$ is bounded away from zero whenever f'(x) is bounded away from zero in the following sense. If for $\{x_{n_k}\} \subset S$, $||f'(x_{n_k})||_1^1 > \varepsilon$ for some $\varepsilon > 0$; then there exists a subsequence of $\{x_{n_k}\}$ and $\delta > 0$ such that $[f'(x_{n_k}), \phi(x_{n_k})] \geq \delta$.

In the theorem which follows, $-\phi(x)$ serves to define the descent direction.

AFFEL-TE-72-77

2. BASIC TERDERA ON GENERALIZED DESCENT ALGORITHMS

Many of the ideas in the following theorem, were stimulated by two papers by A. A. Goldstein (References 6, 7). In particular the definition of ¢ given above and the form of the conclusions of the theorem are identical to those of Goldstein. The hypotheses of the theorem are changed to specialize to the case of a single-dimensional search at each iteration. Thus the proof of part (a) of the theorem is changed. The proofs of parts (b) and (c) follow Goldstein. The following additional remark on the hypotheses is in order.

In the current setting where X is a normed linear space, the assumption that f' is uniformly continuous on S may be replaced with the equivalent conditions that f is uniformly differentiable and that f' is bounded on S (Reference $1\frac{1}{2}$, p, $\frac{1}{2}$ 5).

Theorem I

Assume S is bounded in X, f is bounded below on S and the Frechetderivative f' of f exists, is uniformly continuous on S and bounded on S.

Set $x_{n+1} = x_n$ when $[f^*(x_n), \phi(x_n)] = 0$, otherwise choose $\rho = \rho_n$ to minimize $\{f(x_n - \rho\phi(x_n)): \rho: 0\}$ and set $x_{n+1} = x_n - \rho_n\phi(x_n)$. Then

- (a) $f(x_n) + L$, $f'(x_n) + 0$,
- (b) if $\{x_n\}$ has cluster points, every cluster point z satisfies f(z)=1, f'(z)=0.
- (c) If f' has finitely many zeros on S, S is compact, and $||x_{n+1} x_n|| \to 0$, then the sequence $\{x_n\}$ converges.

Proof

(a) For x(S and f'(x) \neq 0, [f'(x), φ (x)]>0. In general, since f is differentiable,

AFFEL-Th-72-77

$$[f'(x), h] = \lim_{t\to 0} \frac{1}{t} (f(x+th) - f(x))$$

$$= \lim_{t\to 0^{-}} \frac{1}{t} (f(x+th) - f(x))$$

$$= \lim_{t\to 0^{+}} \frac{1}{-p} (f(x-ph) - f(x)),$$

$$= \lim_{t\to 0^{+}} \frac{1}{-p} (f(x-ph) - f(x)),$$

so in particular,

$$0 < [f'(x), \phi(x)] = \lim_{\rho \to 0^{+}} \frac{1}{\rho} (f(x-\rho\phi(x)) - f(x)).$$

By the uniform differentiability of f there exists a $\rho_0>0$, independent of x, such that

$$0 \leftarrow \frac{1}{-\rho_{o}} (f(x-\rho_{o}\phi(x)) - f(x))$$

50,

$$0 \cdot f(x - \rho_0 \phi(x)) - f(x)$$

or

$$f(x-\rho_0\phi(x))< f(x)$$
.

In particular for $x = x_n(S, f(x_n - \rho_0\phi(x_n)) < f(x_n) \le f(x_0)$.

Also for $f'(x) \neq 0$, $[f'(x), \phi(x)] \neq 0$, hence $||\phi(x)|| > 0$. Since S is assumed bounded, there exists a $\rho_a > 0$ such that for all $\rho > \rho_a$, $x - \rho \phi(x) \not\in S$. For on the contrary assumption, for every N, no matter how large, there exists a $\rho_N > N$ such that $y = x - \rho_N \phi(x) \in S$. Then $||y-x|| = \rho_N ||\phi(x)|| > N ||\phi(x)||$ is unbounded which contradicts the boundedness of S.

Now for $x_n \notin S$ such that $f^*(x_n) \neq 0$, $F(\rho) = f(x_n - \rho \phi(x_n)) - f(x_n)$ has a minimum at $\rho_n \notin [0, \rho_a]$ since F is continuous and $[0, \rho_a]$ is compact. Furthermore, $\rho_n \neq 0$ since $F(\rho_0) = f(x_n - \rho_0 \phi(x_n)) - f(x_n) < F(0)$. Thus the sequences $\{\rho_n\}$ and $\{x_n\}$ are well defined and $f(x_{n+1}) = f(x_n - \rho_n \phi(x_n)) < f(x_n)$ which implies the sequence $\{x_n\} \subseteq S$, and the sequence

 $\{f(x_n)\}\$ is strictly decreasing. Since f is assumed bounded below $f(x_n) + L.$

To show $f'(x_n) \neq 0$, we suppose the contrary, them there exists a subsequence $\{x_{n_k}\} \subseteq \{x_n\}$ such that $f'(x_{n_k})$ is bounded away from zero, which implies that a subsequence of $[f'(x_{n_k}), \phi(x_{n_k})]$ is bounded away from zero. Without loss of generality denote this subsequence by $\{x_{n_k}\}$. Then there exists an $\epsilon_0 > 0$ such that $\epsilon_0 < [f'(x_{n_k}), \phi(x_{n_k})] = \lim_{\rho \neq 0^+} \frac{1}{-\rho} (f(x_{n_k} - \rho \phi(x_{n_k})) - f(x_{n_k}))$. Since f is uniformly differentiable, there exists a $\rho_b > 0$, independent of x (i.e., x_{n_k}), such that

$$\frac{\epsilon_0}{2} < \frac{1}{-\rho_b} (f(x_{n_k} - \rho_b \phi(x_{n_k})) - f(x_{n_k})),$$

$$\frac{-\rho_b \varepsilon_0}{2} > f(x_{n_k} - \rho_b \phi(x_{n_k})) - f(x_{n_k}),$$

or
$$f(x_{n_k} - \rho_b \phi(x_{n_k})) < f(x_{n_k}) - \frac{\rho_b \epsilon_0}{2}$$
.

Then
$$f(x_{n_k+1}) \le f(x_{n_k} - \rho_b \phi(x_{n_k})) < f(x_{n_k}) - \rho_b \frac{\epsilon_0}{2}$$
.

Since $f(x_n) + L$,

$$f(x_{n_k}) < f(x_{n_{k-1}}) < \cdots < f(x_{n_{k-1}+1}) < f(x_{n_{k-1}}) - \frac{\rho_b \epsilon_o}{2}$$

tat is,

$$f(x_{n_k}) < f(x_{n_{k-1}}) - \frac{\rho_b \varepsilon_o}{2},$$

$$< f(x_{n_{k-2}}) - 2 \frac{\rho_b \varepsilon_o}{2},$$

$$< f(x_{n_{k-3}}) - 3 \frac{\rho_b \varepsilon_o}{2},$$

$$...$$

$$< f(x_{n_o}) - k \frac{\rho_b \varepsilon_o}{2}.$$

This contradicts the assumption that f is bounded below, hence $f'(x_n)+0$.

AFFEL-78-72-77

- (b) If z is a cluster point of $\{x_n\}$ there exists a subsequence $\{x_{n_k}\}\subseteq \{x_n\}$ such that $x_{n_k}\to z$ (in norm). Since f and f' are continuous, $f(x_n)\to 1$ and $f'(x_n)\to 0$, it follows that f(z)=1 and f'(z)=0.
- (c) Since f'(z) = 0 for every cluster point z of $\{x_n\}$, the number of roots of f' on S is equal to or greater than the number of cluster points of $\{z_n\}$. If f' has a unique root z on S then $\{x_n\}$ converges to it; for otherwise, since S is assumed compact, there exists at least one cluster point z_1 of $\{x_n\}$ in S where $f'(z_1) = 0$ (by (b)). If $z_1 \neq z$, the root of f' is not unique. If the number of roots of f' on S is finite then the number of cluster points of $\{x_n\}$ is finite also.

Let z_i , i = 1,2,...,k be the cluster points of $\{x_n\}$, let $\varepsilon = \min\{||z_i - z_j||: i \neq j, i, j = 1,2,...k\}$, let $S(z_i,\varepsilon/3)$ denote the open sphere of radius $\varepsilon/3$ centered at z_i . Since S is assumed compact k the set $\{x_n\} - \bigcup_{i=1}^{k} S(z_i,\varepsilon/3)$ contains a finite number of points, say m.

Since $||\mathbf{x}_{n+1} - \mathbf{x}_n|| \to 0$, by assumption, there exists an N such that $||\mathbf{x}_{p+1} - \mathbf{x}_p|| < \varepsilon/3n$ for all p>N.

Now, since the z_i , i=1,...k are cluster points of $\{x_n\}$, there are members of $\{x_n\}$ in each $S(z_i,\varepsilon/3)$ for which w>N. Therefore fix n>N such that $x_n \in S(z_i,\varepsilon/3)$, for some fixed i, and $x_{n+1} \notin S(z_i,\varepsilon/3)$. Let x_q be the next member of the sequence $\{x_n\}$ (i.e., q>n) such that $x_q \in S(z_j,\varepsilon/3)$, for some $j \neq i$ and $x_{q-1} \notin S(z_j,\varepsilon/3)$.

Since $x_n \in S(z_i, \varepsilon/3)$ and $x_q \in S(z_j, \varepsilon, 3)$, $i \neq j$, $||x_n - x_j|| > \varepsilon/3$. On the other hand, since n, q > N, we have

$$\begin{aligned} ||\mathbf{x}_{q} - \mathbf{x}_{n}|| &\leq ||\mathbf{x}_{q} - \mathbf{x}_{q-1}|| + ||\mathbf{x}_{q-1} - \mathbf{x}_{q-2}|| + \cdots \\ &+ ||\mathbf{x}_{n+2} - \mathbf{x}_{n+1}|| + ||\mathbf{x}_{n+1} - \mathbf{x}_{n}|| \\ &= \sum_{p=n}^{q-1} ||\mathbf{x}_{p+1} - \mathbf{x}_{p}|| < \sum_{p=n}^{q-1} \frac{\varepsilon}{3m} = (q-n) \frac{\varepsilon}{3m} .\end{aligned}$$

AFFOL-IR-TZ-TT

Then $\frac{\varepsilon}{3} < ||x_q - x_n|| < (q-n) \frac{\varepsilon}{3n}$ or n < q-n.

Now the set of points x_{n+1} , x_{n+2} , ..., x_{q-1} belong to the set $\{x_n\}$ - $\frac{k}{i=1}$ $S(z_i, \epsilon/3)$. But there are (q-1) - (n+i) + i = q-n-1 points in the set $\{x_{n+1}, \dots, x_{q-1}\}$, and by supposition n points in the set $\{x_n\}$ - $\frac{k}{i=1}$ $S(z_i, \epsilon/3)$. Thus $m \ge q-n-1$. On the other hand we have just shown $q-n \ge n$, thus $m \ge q-n-1 \ge n-1$, that is to say q-n-1 = n. Heuristically, all the points of the set $\{x_n\}$ - $\bigcup_{i=1}^{k}$ $S(z_i, \epsilon/3)$ have been accounted for, or "used up". But by the same argument, there exists an $n \ge q$ and a $q^i \ge n^i$ such that $x_n \in S(z_j, \epsilon/3)$ and $x_{n+1} \notin S(z_j, \epsilon/3)$, for the j defined by q above (i.e., for which $x_q \in S(z_j, \epsilon/3)$), and $x_{n+1} \notin S(z_i, \epsilon/3)$, $1 \ne j$, and $1 \ne j$. Then there are again $1 \ne j$ and $1 \ne j$ and $1 \ne j$. Then there are again $1 \ne j$ which also belong to the set $1 \ne j$. Since $1 \ne j$ is $1 \ne j$. There now are $2 \ne j$ points in this set which by supposition contained only $1 \ne j$. This contradiction persists unless we suppose $1 \ne j$ has a unique cluster point,

Comments on the Basic Theorem

The first conclusion of the theorem, $f(x_n)+L$, can be obtained with the weaker assumption that f is differentiable and without the assumption that ϕ is bounded, in which case both ρ_o and ρ_a are dependent on x (or x_n). However, both uniform differentiability of f, which follows from the uniform continuity of f, and boundedness of ϕ were used to assure the existence of a ρ_b which uniformly bounds the differential away from zero in the proof leading to $f'(x_n)+0$.

AFFEL-118-12-17

If X is finite dimensional then S being closed and bounded is necessarily compact and $\{x_n\}\subseteq S$ has cluster points. Hence conclusion (b) of the theorem applies.

3. APPLICATION OF THE MASIC THRORM TO COMON DESCRIT ALCORITINS

In the applications which follow, the boundedness of ϕ will be derived from the boundness of the Frechet derivative f' of f on S. In particular, assume now that X is a Hilbert space so that f'(x) can be represented by its gradient $\nabla f(x)$ in X. Then also since in a Hilbert space every bounded set is weakly compact, the boundedness of f' on S follows from its continuity on S (Reference .3, p. 19).

Corollary I

Let Q be a positive definite continuous linear operator on X, let $\phi(x) = Q^{2}f(x)$, then Theorem I applies.

Proof:

Since f has the required properties, all that remains is to show ϕ is bounded and satisfies conditions (i) and (ii).

By assumption f' is uniformly continuous and bounded on S. This together with Q being a continuous linear operator on X implies \$ = QVf is bounded on S.

Q positive definite implies there exists a m>0 such that $\|z\|^2 \le [z, Qz]$ for all z in X,

hence

The stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of the stransformation of

 $[f'(x), \phi(x)] = [f'(x), Qf'(x)] \ge m ||f'(x)||^2 \ge 0$ and condition (i) is satisfied.

To show condition (ii) assume an $\varepsilon>0$ and choose $\delta=m\varepsilon^2$. Then $[f'(x), \phi(x)] < \delta=m\varepsilon^2$ implies $m | |\nabla f(x)| |^2 \le |\nabla f(x), \nabla f(x)| < \delta=m\varepsilon^2$ or $|\nabla f(x)| | < \varepsilon$.

AFFEL-19-72-77

In order to demonstrate that the theory developed thus far is a generalization of descent algorithms, the Theorem and Corollary will be applied to two common algorithms; the gradient and the Fletcher-Powell methods.

Application to Gradient Algorithms

Let Q in Corollary I be the identity operator, i.e., $\phi(x) = \nabla f(x)$, this is the usual gradient algorithm.

The operator Q need not remain fixed so long as it is uniformly positive definite on S. That is to say there exists a constant m>0, independent of x, such that $[f^1(x), \phi(x)] = [f^1(x), \phi(x)] \nabla f(x) | \nabla f(x)|^2$ for all x in S.

Application to Fletcher-Powell

Let X be finite dimensional, let $Q = Q(k) = H_k$ as defined by R. Fletcher and M.J.D. Powell in Reference 4. Using the algorithm of Fletcher-Powell it is easy to prove H_k remains positive definite at each iteration unless the algorithm terminates in a finite number of steps. Then $\phi(x_n) = H_n \nabla f(x_n)$ satisfies condition (1): $[f'(x_n), \phi(x_n)] \ge 0$. Now if the H_k are uniformly positive definite, or if there exist constants m>0 and p>0, independent of k, such that $[x, H_k x] \ge m||x||^p$, then condition (ii) is satisfied. For given an $\epsilon > 0$, choose $\delta = \max^p$, then $m||f'(x)||^p \le [f'(x), H^k f'(x)] < \delta = m\epsilon^p$ implies $||f'(x)|| < \epsilon$. Then the theorem applies to the Fletcher-Powell algorithm. Furthermore, since S is closed and bounded, it is compact, and the sequence $\{x_k\}$ necessarily has cluster points so (b) of the theorem applies ithout further assumption.

In the above application to Flercher-Powell's method and in the following corollary, the function ϕ is not uniquely defined as a

function from S into X. Rather, ϕ is defined only on the sequence $\{x_n\}$ and its value at an arbitrary x in S depends upon which sequence x belongs to. However, the properties of ϕ used in the proof of Theorem I did not depend on ϕ being defined anywhere except at points of the sequence $\{x_n\}$. That is to say the conditions (i) and (ii) need hold only at the points of the sequence $\{x_n\}$. Within this context the search direction will continue to be denoted by $-\phi(x_n)$.

Corollary II

Let
$$\phi(x_0) = \nabla f(x_0)$$
 and

$$\phi(x_n) = Vf(x_n) + K_n \phi(x_{n-1}) \text{ for } n = 1,2...$$

where |Kn| < r<1, then Theorem I applies.

Proof:

Clearly conditions (i) and (ii) are satisfied at x_0 for $[f'(x_0), \phi(x_0)] = ||\nabla f(x_0)||^2 \ge 0$, and given $\epsilon > 0$ choose $\delta = \epsilon^2$.

Claim $[f'(x_{n+1}), \phi(x_n)] = 0$ for n = 0,1,... Since ρ_n is chosen to minimize $\{F(\rho) = f(x_n - \rho\phi(x_n)) - f(x_n): \rho > 0\}, F'(\rho_n) = 0$. But $F'(\rho) = -[f'(x_n - \rho\phi(x_n)), \phi(x_n)]$ and at $\rho = \rho_n$, $x_n - \rho_n \phi(x_n) = x_{n+1}$. Hence $F'(\rho_n) = -[f'(x_{n+1}), \phi(x_n)] = 0$.

Then for $n = 1, 2, \ldots$

$$\begin{split} [f'(x_n), \phi(x_n)] &= [f'(x_n), \nabla f(x_n)] + K_n[f'(x_n), \phi(x_{n-1})] \\ &= [\nabla f(x_n), \nabla f(x_n)] \\ &= ||\nabla f(x_n)||^2 > 0. \end{split}$$

Hence again, given $\varepsilon>0$ choose $\delta=\varepsilon^2$ then $[f'(x_n), \phi(x_n)] < \delta$ implies $||\nabla f(x_n)|| = ||f'(x_n)|| < \varepsilon$. Thus conditions (i) and (ii) are satisfied.

To show boundedness of the $\phi(x_n)$ consider

$$||\phi(x_n)||^2 = [\phi(x_n), \phi(x_n)] = [\nabla f(x_n), \phi(x_n)] + K_n[\phi(x_{n-1}), \phi(x_n)]$$

$$= ||\nabla f(x_n)||^2 + K_n[\phi(x_{n-1}), \phi(x_n)], \text{ and}$$

$$[\phi(x_n), \phi(x_{n-1})] = [\nabla f(x_n), \phi(x_{n-1})] + K_n[\phi(x_{n-1}), \phi(x_{n-1})]$$

$$= K_n ||\phi(x_{n-1})||^2.$$

Thus,
$$||\phi(x_n)||^2 = ||\nabla f(x_n)||^2 + K_n^2 ||\phi(x_{n-1})||^2$$
.

Applying this relation recursively yields

$$\begin{aligned} ||\phi(\mathbf{x}_n)||^2 &= ||\nabla f(\mathbf{x}_n)||^2 + K_n^2 \{||\nabla f(\mathbf{x}_{n-1})||^2 + K_{n-1}^2||\phi(\mathbf{x}_{n-2})||^2\} \\ &= ||\nabla f(\mathbf{x}_n)||^2 + K_n^2 ||\nabla f(\mathbf{x}_{n-1})||^2 + K_n^2 K_{n-1}^2 ||\phi(\mathbf{x}_{n-2})||^2 \end{aligned}$$

$$\begin{aligned} ||\phi(\mathbf{x}_{n})||^{2} &= ||\nabla f(\mathbf{x}_{n})||^{2} + K_{n}^{2}||^{2} f(\mathbf{x}_{n-1})||^{2} + K_{n}^{2} K_{n-1}^{2}||\nabla f(\mathbf{x}_{n-2})||^{2} \\ &+ \cdots + K_{n}^{2} K_{n-1}^{2} \cdots K_{2}^{2}||\nabla f(\mathbf{x}_{1})||^{2} \\ &+ K_{n}^{2} K_{n-1}^{2} \cdots K_{2}^{2} K_{1}^{2}||\nabla f(\mathbf{x}_{0})||^{2}. \end{aligned}$$

Since f'(x) is bounded for $x \in S$ let $M = \sup_{x \in S} ||\nabla f(x)||$

then

$$||\phi(\mathbf{x}_n)||^2 \leq \{1 + K_n^2 + K_n^2 K_{n-1}^2 + \cdots + K_n^2 K_{n-1}^2 \cdots K_2^2 K_1^2\} M^2.$$

Also $K_n^2 \leq r < 1$, so the series

$$\{1 + K_n^2 + K_n^2 K_{n-1}^2 \div \cdots + K_n^2 K_{n-1}^2 \cdots K_2^2 K_1^2\}$$

converges. Hence $||\phi(\mathbf{x}_n)||^2$ is bounded and therefore $\phi(\mathbf{x}_n)$ is bounded.

For a trivial application of Corollary II, K_n may be set to zero for all n, this then generates the usual gradient algorithm. The following application is of much more interest.

Application to Conjugate Gradient

Let
$$K_n = \frac{\left[\nabla f(x_n), \nabla f(x_n)\right]}{\left[\nabla f(x_{n-1}), \nabla f(x_{n-1})\right]}$$
 in Corollary II.

This is the β_{n-1} of Fletcher-Reeves (Reference 5) for X = Rⁿ, and Lasdon, Mitter and Warren (Reference 10), for X a function space.

The condition $|K_n| \leq r < 1$ may appear overly restrictive particularly in the light of the paper by Lasdon, Mitter and Warren which also is set in a Hilbert space H. Unfortunately, the proof of Lasdon, Mitter and Warren contains one minor but significant discrepancy. The proof of their theorem is reproduced and corrected, to illustrate the similarities of the constraints which must be imposed on K_n or β_{n-1} .

Theorem 3 of Lasdon, Mitter and Warren

- If: 1. J(u) is bounded below,
 - 2. J(u) and g(u) = grad J(u) are continuous,
 - 3. $D^2J(u,h,h)$ exists and $|D^2J(u,h,h)| \le m||h||^2$ for m > 0 and all u,h in H,
- 4. $\{u_k\}$ has a cluster point u^* , then the sequence $\{u_k\}$ formed with arbitrary u_0 by applying the conjugate gradient algorithm to J(u) has the following properties:
 - 1. $\lim_{k\to\infty} J(u_k) = J(u^*),$
 - 2. $\lim_{k\to\infty} g(u_k) = g(u^*) = 0$.

Before presenting the proof the following three remarks are pertinent:

- 1. Here $D^2J(u,h,k)$ is the second differential of J at u.
- 2. The form of the algorithm is exactly as given in Section I.2 except here the independent variable is u instead of x.
- 3. From the above proof of Corollary II it follows that $[g_k,s_k] = -||g_k||^2 \text{ and } [s_k,g_{k+1}] = 0 \text{ for all } k.$

Corrected Proof of Theorem 3 of Lasdon, Nitter and Warren
From Taylor's Theorem,

$$J(u_k + \alpha s_k) = J(u_k) + \alpha[g_k, s_k] + 1/2\alpha^2 D^2 J(\xi_k(\alpha), s_k, s_k)$$

AFFEL-TE-72-77

where $\xi(\alpha)$ belongs to the line segment joining u_k and $u_k^{+\alpha s_k}$. Then, using $[s_k,g_k]=-||g_k||^2$ and assumption 3,

$$J(u_k + \alpha s_k) \le J(u_k) - \alpha ||g_k||^2 + 1/2\alpha^2 n ||s_k||^2.$$

Since $\alpha = \alpha_k$ minimizes $J(u_k + \alpha s_k)$,

$$J(u_{k} + \alpha_{k}s_{k}) = J(u_{k+1}) \leq J(u_{k} + \frac{1}{n}s_{k})$$

$$\leq J(u_{k}) - \frac{1}{n} ||g_{k}||^{2} + \frac{1}{2n} ||s_{k}||^{2}.$$

At this point Lasdon, Mitter and Warren assume $||s_k|| = ||g_k||$ in their proof of this theorem. This is clearly inconsistent with $[s_k, g_k] = -||g_k||^2$ for by the Cauchy-Schwarz inequality $|[s_k, g_k]| \le ||s_k|| \, ||g_k||$ where equality holds if and only if s_k is a multiple of g_k . Thus $||g_k||^2 = |[s_k, g_k]| \le ||s_k|| \, ||g_k||$ or $||g_k|| \le ||s_k||$. Specifically, equality holds only when k=0 (i.e., $s_0 = -g_0$) or $g_k = 0$ in which case the solution has been achieved.

Proceeding more carefully, observe that $s_k = -g_k + \beta_{k-1}s_{k-1}$,

$$||s_k||^2 = [s_k, s_k] = -[s_k, g_k] + \beta_{k-1}[s_k, s_{k-1}]$$

$$= ||g_k||^2 + \beta_{k-1}[s_k, s_{k-1}]$$

and
$$[s_{k-1}, s_k] = -[s_{k-1}, g_k] + \beta_{k-1}[s_{k-1}, s_{k-1}]$$

Then

$$\begin{split} \mathbf{J}(\mathbf{u}_{k+1}) & \leq \mathbf{J}(\mathbf{u}_{k}) - \frac{1}{m} ||\mathbf{g}_{k}||^{2} + \frac{1}{2m} (||\mathbf{g}_{k}||^{2} + \beta_{k-1}^{2} ||\mathbf{s}_{k-1}||^{2}) \\ & = \mathbf{J}(\mathbf{u}_{k}) - \frac{1}{2m} ||\mathbf{g}_{k}||^{2} + \frac{1}{2m} \beta_{k-1}^{2} ||\mathbf{s}_{k-1}||^{2}, \end{split}$$

Likewise,

$$J(u_k) \le J(u_{k-1}) - \frac{1}{2m} ||g_{k-1}||^2 + \frac{1}{2m} g_{k-1}^2 ||s_{k-2}||^2,$$

hence

$$\begin{split} J(u_{k+1}) & \leq J(u_{k-1}) - \frac{1}{2n} ||g_{k-1}||^2 + \frac{1}{2n} \beta_{k-2}^2 ||s_{k-2}||^2 \\ & - \frac{1}{2n} ||g_k||^2 + \frac{1}{2n} \beta_{k-1}^2 ||s_{k-1}||^2, \end{split}$$

or in general,

$$J(u_{k+1}) \leq J(u_{j}) - \sum_{i=1}^{k} \frac{1}{2n} ||g_{1}||^{2} + \sum_{i=1-1}^{k-1} \frac{1}{2n} g_{i}^{2} ||s_{i}||^{2},$$

or

$$J(u_{k}) \leq J(u_{0}) + \sum_{i=0}^{k-1} \frac{\beta_{1}^{2}}{2n} ||s_{1}||^{2} - \sum_{i=0}^{k} \frac{1}{2n} ||g_{1}||^{2}$$

$$= J(u_{0}) - \sum_{i=1}^{k} z_{i} - \frac{1}{2n} ||g_{0}||^{2},$$

where
$$z_i = \frac{1}{2n} ||g|!|^2 - \frac{\beta_{i-1}^2}{2n} ||s_{i-1}||^2$$
.

Since J(u) is bounded below $\lim_{k\to\infty} \sum_{i=1}^{k} z_i$ exists and is finite, that is to

say the series $\sum z_i$ converges.

Since $\sum z_i = \sum \left[\frac{1}{2m}||g_i||^2 - \frac{\beta_{i-1}^2}{2m}||s_{i-1}||^2\right]$ it follows that (Reference 13) if the series $\sum z_i$ is absolutely convergent, then each series $\sum \frac{1}{2m}||g_i||^2$ and $\sum \frac{\beta_{i-1}^2}{2m}||s_{i-1}||^2$ is convergent, which implies $||g_i|| + 0$ and $\beta_{i-1}^2||s_{i-1}||^2 + 0$. It also follows that $||s_i||^2 = ||g_i||^2 + \beta_{i-1}^2 ||s_{i-1}||^2 + 0$ from which it follows that $||s_i||$ is bounded. On the other hand if $\sum z_i$ is conditionally convergent (convergent but not absolutely convergent) then each series $\sum \frac{1}{2m}||g_i||^2$ and $\sum \frac{\beta_{i-1}^2}{2m}||s_{i-1}||^2$ is divergent. Of course it still may happen that $||g_i|| + C$ and $\beta_{i-1}^2||s_{i-1}||^2 + 0$ (and hence s_i is bounded).

One way to assure $\sum_{i=1}^{\infty} is$ absolutely convergent is to assure $z_{i} \ge 0$, that is $2m \ z_{i} = ||g_{i}||^{2} - \beta_{i-1}^{2}||s_{i-1}||^{2} \ge 0$. But $||s_{i}||^{2} = ||g_{i}||^{2} + \beta_{i-1}^{2}||s_{i-1}||^{2}$, hence $2m \ z_{i} = 2||g_{i}||^{2} - ||s_{i}||^{2} \ge 0$

requires
$$||s_i||^2 \le 2||s_i||^2$$
 (recall $||s_i||^2 \le ||s_i||^2$). Now $||s_o||^2 = ||s_o||^2$, $||s_o||^2 = ||s_o||^2$, $||s_o||^2 = ||s_o||^2 + \beta_o^2||s_o||^2$

$$= ||s_o||^2 \cdot (1 + \beta_o \frac{||s_o||^2}{||s_o||^2})$$

$$= ||s_o||^2 \cdot (1 + \beta_o)$$
since $\beta_o = ||s_o||^2 / ||s_o||^2$.

Similarly,

$$||s_2||^2 = ||g_2||^2 (1 + \beta_1 ||s_1||^2)$$

$$= ||g_2||^2 (1 + \beta_1 (1 + \beta_0))$$

$$= ||g_2||^2 (1 + \beta_1 + \beta_1 \beta_0)$$

or in general

$$||s_k||^2 = ||g_k||^2 (1 + \beta_k + \beta_k \beta_{k-1} + \cdots + \beta_k \beta_{k-1} \cdots \beta_0)$$

Now if $\beta_i \le r < 1/2$ for all i,

$$1 + \beta_k + \beta_k \beta_{k-1} + \cdots + \beta_k \beta_{k-1} \cdots \beta_0 \le \frac{1}{1-r} < \frac{1}{1-1/2} = 2$$

Then $||s_k||^2 < 2||g_k||^2$, and $z_k > 0$, and absolute convergence of $\sum z_k$ follows from its convergence. The condition $\beta_i \le r < 1/2$ is stronger than the condition $|K_n| \le r < 1$ assumed in Corollary II.

Application to Conjugate Gradient in RN

The classic proof of the conjugate gradient algorithm for function minimization in finite dimensional spaces (R^N) is based on the Gram-Schmidt Orthogonalization procedure. Originally the method was developed for the solution of systems of linear equations. The extension to the problem of minimizing a quadratic function on R^N is well-known. In theory, the conjugate gradient method finds the minimum of a quadratic function in at most N steps. However in practice, be-

AFFEL-TR-72-TT

cause of round-off error, N+1 steps are used to obtain the "exact" solution.

Application of the conjugate gradient algorithm to the problem of minimizing an arbitrary function in \mathbb{R}^K usually follows the same procedure, viz., conjugate directions are used for M+1 steps and if the minimum has not been obtained, the algorithm is re-started with the last best estimate (\mathbf{x}_{M+1}) . In terms of Corollary II, this means $K_{M+1} = K_{2M+2} = \cdots = 0$. This condition may be used to impose a somewhat weaker but less instructive condition on the K_{n+1} in Corollary II, specifically

$$(1 - K_N^2 + K_N^2 K_{N-1}^2 + \cdots + K_N^2 K_{N-1}^2 \cdots K_1^2) \le A$$

 $(1 + K_{2N+1}^2 + K_{2N+1}^2 K_{2N}^2 + \cdots + K_{2N+1}^2 K_{2N}^2 \cdots K_{N+2}^2) \leq A \text{ etc. for }$ each sub-cycle. Then $||\phi(x_k)||^2 \leq A M^2$ for all k where $||\nabla f(x_k)||^2 \leq M^2$.

As a final remark it should be noted that in Corollary II the condition $|K_n| \le r < 1$ need not be satisfied for all n but only for all n beyond some point. In the conjugate gradient algorithm $K_n = \frac{||\nabla f(x_n)||^2}{||\nabla f(x_{n-1})||^2} \le r < 1 \text{ implies the gradient not only converges to}$ zero but $||\nabla f(x_n)||$ forms (eventually) a strictly monotone decreasing sequence.

SCHOOL HIE

A PARK-ONE MEDEOD OF FUNCTION MENT FEXAULOR

A particular property possessed by both the method of conjugate gradients and the method of Fletcher-Powell is that either method obtains the minimum of a positive definite quadratic form; viz., minimizes $f(x) = f_0 + a^t x + 1/2 x^t Gx$, $x \in \mathbb{R}^n$, in a finite number of steps excepting round-off errors. This is accomplished in the conjugate gradient algorithm by means of the Gram-Schmidt orthogonalization process (Reference 1). In the method of Fletcher-Powell, this is accomplished by generating the inverse of the matrix G, specifically, $H_{ij} = G^{-1}$. In this section another method for generating G^{-1} is presented which does not require a single-dimensional search for a minimum as do the methods of Fletcher-Powell and conjugate gradients.

Although this author arrived at the method independently, the algorithm is essentially the same as Davidon's variance algorithm (Reference 3). Although Davidon's proofs are valid, they provide little insight for the user on how the algorithm is structured or why the method works.

The derivation presented here clarifies the structure of the algorithm by emphasizing how the structure leads to the desired properties at each iteration. Finally, a minor change is incorporated into Davidon's algorithm which circumvents one difficulty which may be encountered when the algorithm is applied to a computational problem.

1. BASIC RANK-ONE METHOD

In Reference 15, Herbert S. Wilf presents a method for matrix inversion based on the equation $(A + uv')^{-1} = A^{-1} - \frac{A^{-1}uv'A^{-1}}{1 + v'A^{-1}u}$

where the prime (') denotes matrix (vector) transpose and it is assumed \mathbf{A}^{-1} exists. This equation is easily verified by computing

13501-18-12-FT

$$(A + wv^{i})(A^{-1} - \frac{A^{-1}wv^{i}A^{-1}}{1 + v^{i}A^{-1}u}) = I.$$

Note that A + wv differs from A by a matrix, wv, of rank one. Therefore, the above inversion technique and the algorithm developed below might well be called Rank-One Methods.

D_n denote the partial sum inverses

$$D_n = C_n^{-1} = (C_0 + \sum_{i=1}^n u_i v_i^i)^{-1}.$$

Then $C_{n+1} = C_n + u_{n+1}v_{n+1}^t$

and
$$D_{n+1} = (C_0 + \sum_{i=1}^{n+1} u_i v_i^*)^{-1}$$

$$= (C_n + u_{n+1} v_{n+1}^*)^{-1}$$

$$= (D_n^{-1} + u_{n+1} v_{n+1}^*)^{-1}$$

$$= D_n - \frac{D_n u_{n+1} v_{n+1}^* D_n}{1 + v_{n+1}^* D_n u_{n+1}}.$$

It follows that $D_N = C_N^{-1} = (C_0 + \sum_{i=1}^{N} u_i v_i)^{-1} = C^{-1}$ providing all computa-

tions may be carried out. Obviously the procedure fails if the denominator term $1+v_{n+1}^{\dagger}D_{n}u_{n+1}$ is zero. The implication of such a condition can be deduced from the results of the following Lemma.

Lemma I

The states of the contraction of the state o

Assume A is nonsingular, then the matrix B = A + uv' is singular if and only if $1 + v'A^{-1}u = 0$.

Proof

Necessity: Assume B is singular, then there exists an $x \neq 0$ such that $Bx = Ax + uv^{\dagger}x = 0$. Then $A^{-1}Bx = x + A^{-1}uv^{\dagger}x = 0$ and

AFFEL-TB-TZ-TT

 $v^*A^{-1}Bx = v^*x (1 + v^*A^{-1}w) = 0$. Now $v^*x \neq 0$, for them Bx = Ax = 0but A is nonsingular. Therefore, $(1 + v^*Aw) = 0$.

Selficiency: Assume $1 + v^*Au = 0$ and set $x = A^{-1}u \neq 0$ for otherwise u = 0 and $v^*A^{-1}u = 0$ which contradicts $1 + v^*A^{-1}u = 0$. Then $B_X = Ax + uv^*x$ $= u + uv^*A^{-1}u$ $= (1 + v^*A^{-1}u) u$ = 0

2. STENCTURE OF THE BANK-ONE MINIMIZATION ALGORITHM

An algorithm will be constructed for minimizing the quadratic $f_{0,m} f(x) = f_0 + a^t x + 1/2 x^t G_{\ell}$, $x \in \mathbb{R}^M$ and G a positive definite symmetric matrix. The method is based on the above technique for generating G^{-1} . Let $B_0 = C_0^{-1}$, where C_0 is an arbitrary positive definite symmetric matrix such as the identity I. Pick an arbitrary x and let $g_0 = a + Gx_0 = \text{grad } f(x_0)$. Assume $g_0 \neq 0$ for otherwise x_0 is the required solution. Set $\sigma_0 = -D_0g_0$, $x_1 = x_0 + \sigma_0$, $g_1 = \text{grad } f(x_1) = a + Gx_1$ and $y_0 = g_1 - g_0$. Then $g_1 = a + G(x_0 + \sigma_0) = g_0 + G\sigma_0$ and $y_0 = G\sigma_0$. An improved estimate c_1 to c_1 is sought such that $c_1 \sigma_0 = c_0 = c_0$, where C_1 has the form $C_1 = C_0 + u_1v_1'$. Since G is symmetric, being the Hessian of f, the added constraint $u_1 = v_1$ is imposed to assure symmetry of c_1 . Then the condition to be satisfied is $c_1\sigma_0 = c_0\sigma_0 +$ $u_1 u_1^{\dagger} \sigma_0 = y_0$, or $u_1 = \frac{y_0 - C_0 \sigma_0}{u_1^{\dagger} \sigma_0}$. But $-C_0 \sigma_0 = -C_0 (-D_0 g_0) = g_0$, so $u_1 = \frac{g_1}{u_1^* \sigma_0} = kg_1$. The value of the scalar k is easily determined from $kg_1 = \frac{g_1}{u_1^* \sigma_0}$ to be $k^2 = \frac{1}{g_1^* \sigma_0}$ and C_1 takes the very simple form $c_1 = c_0 + \frac{g_1 g_1^*}{g_1^* \sigma}$. Applying the rank-one inversion formula to c_1 yields

17701-TB-72-77

$$B_{1} = c_{1}^{-1} = c_{0}^{-1} - \frac{c_{0}^{-1} \frac{s_{1}s_{1}^{*}}{s_{1}^{*}c_{0}} c_{0}^{-1}}{1 + \frac{s_{1}^{*}c_{0}^{-1}s_{1}}{s_{1}c_{0}}}, \text{ or } B_{1} = B_{0} - \frac{B_{2}s_{1}^{*}B_{0}}{s_{1}^{*}c_{0}^{-1} + s_{1}^{*}B_{0}s_{1}}.$$

But
$$s_1^* D_0 s_1 = s_1^* D_0 (3. + y_0)$$

= $s_1^* D_0 y_0 - s_1^* c_0$

so
$$D_1 = D_0 - \frac{D_0 E_1 E_1^{*} D_0}{E_1^{*} D_0 y_0}$$
.

Just as
$$C_1 \sigma_0 = y_0$$
, $D_1 y_0 = \sigma_0$ for $D_1 y_0 = D_0 y_0 - \frac{D_0 g_1 g_1^* D_0 y_0}{g_1^* D_0 y_0}$

$$= D_o(g_1 - S_o) - D_cg_1 = -D_og_o = \sigma_o$$

or simply by applying $D_1 = C_1^{-1}$ to $C_1\sigma_0 = y_0$.

Since D_1 is the current best estimate of G^{-1} , it is reasonable to continue the process by computing $\sigma_1 = -D_1g_1$, $x_2 = x_1 + \sigma_1$, $g_2 = \operatorname{grad} f(x_2)$ and $y_1 = g_2 - g_1$ and considering an improved estimate C_2 to G obtained from C_1 , g_2 and y_1 in precisely the same manner. More generally suppose x_n, D_n, g_n and C_n have been obtained. Set

$$\sigma_{n} = -D_{n}g_{n},$$

$$x_{n+1} = x_{n} + \sigma_{n},$$

$$g_{n+1} = \text{grad } f(x_{n+1}),$$

$$y_{n} = g_{n+1} - g_{n} = G\sigma_{n}.$$

Define
$$D_{n+1} = D_n - \frac{D_n g_{n+1} g_{n+1}^{\dagger} D_n}{g_{n+1}^{\dagger} D_n y_n}$$

and
$$C_{n+1} = C_n + \frac{g_{n+1}g_{n+1}^*}{g_{n+1}^*\sigma_n}$$
.

Clearly
$$C_{n+1}\sigma_n = C_n\sigma_n + g_{n+1} = -g_n + g_{n+1} = y_n = G\sigma_n$$
 and $D_{n+1}y_n = D_ny_n - D_ng_{n+1} = D_ng_{n+1} - D_ng_n - D_ng_{n+1} = \sigma_n = G^{-1}y_n$.

17701-119-72-77

An important property of the algorithm is if $D_m y = \sigma = C^{-1} y$ for an arbitrary y then $V_{m+1}y = D_m y$. This property is demonstrated by showing $g'_{m+1}D_m y = 0$ as follows. Since $g'_{m+1} - g'_m = C(x_{m+1} - x_m)$,

$$\sigma_{m} = x_{m+1} - x_{m} = C^{-1}(g_{m+1} - g_{m}) = -D_{m}g_{m},$$

or
$$-D_{m}g_{m} - C^{-1}(g_{m+1} - g_{m}) = 0$$
.

Then
$$D_m g_{m+1} = D_m g_m - G^{-1}(g_{m+1} - g_m) + D_m g_{m+1}$$

= $(D_m - G^{-1})(g_{m+1} - g_m)$,

and
$$g_{n+1}^{\prime} D_n y = (g_{n+1} - g_n)^{\prime} (I_n - G^{-1}) y = 0.$$

That is to say, if D_n agrees with G^{-1} for some vector y, then so does D_{n+1} .

Clearly D_{n+1} cannot be computed if the term $g_{n+1}^*D_ny_n$ vanishes. If this occurs as above because D_n agrees with G^{-1} on y_n , then $D_ny_n=G^{-1}y_n=G^{-1}(g_{n+1}-g_n)$ $=G^{-1}(a+f_{n+1}-a-Gx_n)$

= $x_{n+1} - x_n = \sigma_n = -D_n g_n$. But also, $D_n y_n = D_n g_{n+1} - D_n g_n$; therefore $D_n g_{n+1} = 0$ and if D_n is non-singular $g_{n+1} = 0$ which in turn implies x_{n+1} is the desired solution.

By Lemma I, D_{n+1} exists so long as the denominator $(1 + v^*A^{-1}u$ in the Lemma) does not vanish. In the case of the D_k^*s , D_{n+1} exists so long as $g_{n+1}^*\sigma_n + g_{n+1}^*D_ng_{n+1} = g_{n+1}^*D_ny_n \neq 0$. Furthermore, C_{n+1} can be obtained from D_{n+1} by applying the same inversion technique, hence the denominator in the expression $C_{n+1} = C_n + \frac{g_{n+1}g_{n+1}^*}{g_{n+1}^*\sigma_n}$ must not vanish.

There are then two requirements for the existence of $D_{n+1} = C_{n+1}^{-1}$ and $C_{n+1} = D_{n+1}^{-1}$, namely: $g_{n+1}^{\dagger} D_n y_n \neq 0$ and $g_{n+1}^{\dagger} \sigma_n \neq 0$.

In the following it is assumed, for the moment, that these requirements are satisfied.

The algorithm may be generalized to one for minimizing an arbitrary function f defined on \mathbb{R}^N as follows:

Initially: choose an arbitrary
$$x_o$$
, set D_o = T and compute g_o = grad $f(x_o)$.

Iteratively: Set $\sigma_n = -D_n g_n$ and $x_{n+1} = x_n + \sigma_n$.

Compute $g_{n+1} = grad \ f(x_{n+1})$,

set $y_n = g_{n+1} - g_n$ and

compute $D_{n+1} = D_n - \frac{D_n g_n}{g_{n+1}^2 D_n f_n}$.

As an algorithm for minimizing an arbitrary function, precautions must be taken to avoid instances where the denominator will vanish. In addition, for an arbitrary function, D_n may not be positive definite and hence $\sigma_n = -D_n g_n$ may not be a descent direction. Therefore, additional precautions must be incorporated into the algorithm to assure applicability to arbitrary functions at the expense of increased complexity.

The recursive relation for D_{n+1} may be rewritten as follows:

$$D_{n+1} = D_n - (\frac{g_{n+1}^{\dagger} D_n g_{n+1}}{g_{n+1}^{\dagger} D_n y_n}) \frac{D_n g_{n+1} g_{n+1}^{\dagger} D_n}{g_{n+1}^{\dagger} D_n g_{n+1}}$$

or
$$D_{n+1} = D_n + (\lambda_n - 1) \frac{D_n g_{n+1} g_{n+1}^{\dagger} D_n}{g_{n+1}^{\dagger} D_n g_{n+1}}$$

where
$$\lambda_n = 1 - \frac{g_{n+1}^{\dagger} D_n g_{n+1}}{g_{n+1}^{\dagger} D_n y_n}$$
.

This form for D_{n+1} is essentially the recursive relation of Davidon's Variance Algorithm except that A_n is chosen to assure both D_{n+1} and D_{n+1}^{-1} (i.e., C_{n+1}) remain positive definite at each iteration, where the

above relation for λ_n is used whenever possible. Specifically, λ_n is chosen such that $2x^*D_nx\le x^*D_{n+1}x\le x^*D_nx$ for all x, where $0<x<1<\beta$.

As Davidon's Variance Algorithm is presented in Reference 3, D_{n+1} (or V^+ in Davidon's notation) is constructed from a test value x^\pm for x_{n+1} . If $f(x^\pm) \geq f(x_n)$, then x_n is taken for x_{n+1} , that is the estimate x^\pm is discarded. As pointed out by Davidon in a footnote (Reference 3, p. 406) the algorithm can become trapped in a loop. Indeed, although the poor estimate x^\pm is not used, the gradient of f at x^\pm is used to modify the estimate, D_n , of G^{-1} . Since for most problems f(x) can be computed much more rapidly than $g(x) = \operatorname{grad} f(x)$, several test computations of f(x), without the corresponding gradient, can be made without undue increase in computation time. The following modification to Davidon's Variance algorithm provides for a search for an improved estimate for x at each iteration, if necessary, using several computations of f(x) before computing g(x) and updating the estimate of G^{-1} . Since this search introduces a major change in the algorithms, the form of the general term is re-derived as follows.

Assume x_n , g_n , $D_n = C_n^{-1}$ have been obtained where $g_n = \operatorname{grad} f(x_n)$, D_n is positive definite, and $f(x) = f_0 + a^*x + 1/2 x^*Gx$. Pick an α_n such that $f(x_n - \alpha_n D_n g_n) < f(x_n)$ using $\alpha_n = 1$ whenever possible, otherwise α_n is reduced (for example $\alpha_n = 1/2$, 1/4,...) until f is decreased. This is possible since $[f'(x_n), -D_n g_n] = -g'_n D_n g_n < 0$. Set $\sigma_n = -\alpha_n D_n g_n$, $x_{n+1} = x_n + \sigma_n$ and compute $g_{n+1} = \operatorname{grad} f(x_n) = g_n + G\sigma_n$ and $y_n = g_{n+1} - g_n = G\sigma_n$. An improved estimate C_{n+1} to G of the form $C_{n+1} = C_n + u_n u_n^*$ is sought such that $C_{n+1}\sigma_n = G\sigma_n = g_n + g_n = g_{n+1} - g_n$. As before, $C_{n+1}\sigma_n = C_n\sigma_n + u_n u_n^*\sigma_n = g_{n+1} - g_n$ from which $u_n = \frac{g_{n+1} - g_n - C_n\sigma_n}{u_n^*\sigma_n}$,

or
$$u_n = \frac{g_{n+1} - (1 - \alpha_n)g_n}{u_n^* \sigma_n}$$

since $C_n \sigma_n = -\alpha_n C_n (D_n g_n) = -\alpha_n g_n$

For clarity, set $v_n = g_{n+\frac{1}{2}} - (1 - c_n)g_n$ and observe u_n has the form

$$u_n = \frac{v_n}{k} . \quad \text{Then } = \frac{v_n}{k} = \frac{v_n}{u_n^* \sigma_n} = \frac{v_n}{v_n^* \sigma_n/k} \quad \text{and } u_n^* \sigma_n = \frac{v_n^* \sigma_n}{k} = \frac{v_n^* \sigma_n}{v_n^* \sigma_n/k} = k,$$

from which $k^2 = v_n^{\dagger} \sigma_n$.

Now C_{n+1} may be written $C_{n+1} = C_n + \frac{\mathbf{v}_n \mathbf{v}_n^t}{\mathbf{v}_n^t \sigma_n}$ and an application of the

rank-one inversion formula yields $D_{n+1} = C_{n+1}^{-1} = D_n - \frac{D_n v_n v_n^* D_n}{v_n^* \sigma_n + v_n^* D_n v_n}$.

Since
$$\mathbf{v}_{n}^{\dagger}\mathbf{D}_{n}\mathbf{v}_{n} = \mathbf{v}_{n}^{\dagger}\mathbf{D}_{n}[\mathbf{g}_{n+1} - \mathbf{g}_{n} + \alpha_{n}\mathbf{g}_{n}]$$

$$= \mathbf{v}_{n}^{\dagger}\mathbf{D}_{n}\mathbf{y}_{n} + \alpha_{n}\mathbf{v}_{n}^{\dagger}\mathbf{D}_{n}\mathbf{g}_{n}$$

$$= \mathbf{v}_{n}^{\dagger}\mathbf{D}_{n}\mathbf{y}_{n} - \mathbf{v}_{n}^{\dagger}\sigma_{n},$$

 \mathbf{D}_{n+1} may be simplified to $\mathbf{D}_{n+1} = \mathbf{D}_n - \frac{\mathbf{D}_n \mathbf{v}_n \mathbf{v}_n^{\dagger} \mathbf{D}_n}{\mathbf{v}_n^{\dagger} \mathbf{D}_n \mathbf{y}_n}$.

This form of the algorithm can be shown to have the same properties as Davidon's form. For example if D_n agrees with G^{-1} on y, i.e., $D_n y = G^{-1} y = \sigma, \text{ then so does } D_{n+1}.$ This is demonstrated, as before, by showing $v_n^* D_n y = 0$ as follows:

$$\begin{split} \mathbf{D_n}\mathbf{v_n} &= \mathbf{D_{r_i}}(\mathbf{g_{n+1}} - \mathbf{g_n} + \alpha_n\mathbf{g_n}) = \mathbf{D_n}(\mathbf{g_{n+1}} - \mathbf{g_n}) - \mathbf{c_n} \\ \text{but } \sigma_n &= \mathbf{x_{n+1}} - \mathbf{x_n} = \mathbf{G^{-1}} \ (\mathbf{g_{n+1}} - \mathbf{g_n}), \text{ so } \mathbf{D_n}\mathbf{v_n} = (\mathbf{D_n} - \mathbf{G^{-1}})(\mathbf{g_{n+1}} - \mathbf{g_n}), \\ \text{and } \mathbf{v_n^i}\mathbf{D_n}\mathbf{y} = (\mathbf{g_{n+1}} - \mathbf{g_n})' \ (\mathbf{D_n} - \mathbf{G^{-1}})\mathbf{y} = \mathbf{0}. \end{split}$$

This form of the algorithm also suffers the difficulty of D_{n+1} not necessarily being positive definite at each iteration for an arbitrary function. This difficulty can be avoided by applying Davidon's method of assuring boundedness. D_{n+1} may be written as

$$D_{n+1} = D_n - (\frac{g_{n+1}^{\dagger}D_ng_{n+1}}{v_n^{\dagger}D_ny_n}) - \frac{D_nv_nv_n^{\dagger}D_n}{g_{n+1}^{\dagger}D_ng_{n+1}}$$

or
$$D_{n+1} = D_n + (\lambda_n - 1) \frac{D_n v_n v_n^* D_n}{g_{n+1}^* D_n g_{n+1}}$$

where
$$\lambda_n = 1 - \frac{g_{n+1}^{\dagger} D_n g_{n+1}}{v_n^{\dagger} D_n y_n}$$
.

Following Davidon, λ_n is chosen by the above relation whenever possible, otherwise such that $\alpha \leq \lambda_n \leq \beta$ where $0 < \alpha < 1 < \beta < \infty$ in order to assure

$$\alpha x' D_{n+1} x \leq x' D_n x \leq \beta x' D_n x$$

for any x.

3. COMPLETE RANK-ONE ALGORITHM FOR FUNCTION MINIMIZATION

The complete algorithm proceeds as follows:

- (1) Initially set $D_0 = I$ and pick α , β such that $0 < \alpha < 1 < \beta$.

 Choose an arbitrary x_0 and compute $f(x_0)$ and $g(x_0) = \text{grad } f(x_0)$.
- (2) Compute $f(x_n \alpha_n D_n g_n)$ for $\alpha_n = 1, 1/2, 1/4, \dots$ until $f(x_n \alpha_n D_n g_n) < f(x_n)$
- (3) Set $\sigma_n = -\alpha_n D_n g_n$ and $x_{n+1} = x_n + \sigma_n$ Compute $g_{n+1} = g(x_{n+1}) = \text{grad } f(x_{n+1})$ and set $y_n = g_{n+1} - g_n$, $v_n = y_n + \alpha_n g_n = g_{n+1} - (1 - \alpha_n) g_n$ and $\gamma_n = \frac{g'_{n+1} D_n g_{n+1}}{v'_n D_n y_n}$
- (4) If $\gamma_n > 1 \alpha$ set $\lambda_n = \alpha$, if $\gamma_n < 1 - \beta$ set $\lambda_n = \beta$, otherwise set $\lambda_n = 1 - \gamma_n$.
- (5) Set $D_{n+1} = D_n + (\lambda_n 1) \frac{D_n v_n v_n^* D_n}{g_{n+1}^* D_n g_{n+1}}$
- (6) Repeat steps (2) through (5) until selected error tolerances(s) is (are) satisfied. One or more of the following error tests may be used:

$$||\mathbf{g}_{n+1}|| < \varepsilon_1, \ \mathbf{g}_{n+1}^{\dagger} \mathbf{D}_n \mathbf{g}_{n+1} < \varepsilon_2,$$

 $\mathbf{f}(\mathbf{x}_n) - \mathbf{f}(\mathbf{x}_{n+1}) < \varepsilon_3, \ \text{and} \ ||\sigma_n|| < \varepsilon_4.$

It is interesting to note that any algorithm which generates recursive estimates to the inverse of the Hessian as do the Fletcher-Powell algorithm, Davidon's Variance algorithm, or the modified Davidon algorithm presented above can be expected to exhibit numerical difficulties whenever the Hessian or its inverse is singular at the minimizing point. The following examples illustrate simple problems with this property.

Consider the problem of minimizing $f(x,y) = x^2 + y^4$. The first and second derivatives are $\nabla f(x,y) = \begin{pmatrix} 2x \\ 4y^3 \end{pmatrix}$ and

$$\nabla^2 f(x,y) = \begin{bmatrix} 2 & 0 \\ 0 & 12y^2 \end{bmatrix}.$$

Clearly the minimum is at (x,y) = (0,0) where $\nabla f(x,y) = (0,0)$ but $\nabla^2 f(0,0) = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$ is singular. Also,

$$\nabla^{2} \mathbf{f}(\mathbf{x}, \mathbf{y})^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & \frac{1}{12\mathbf{y}^{2}} \end{bmatrix} \xrightarrow{\mathbf{y} \to 0} \begin{bmatrix} 1/2 & 0 \\ 0 & \infty \end{bmatrix}.$$

Any algorithm which attempts to estimate $\nabla^2 f(x,y)^{-1}$ can well be expected to have terms which tend to become unbounded as the solution is approached.

For an example of the inverse condition, consider minimizing

 $f(x,y) = x^2 + y^{4/3}$. Then, formally, $\nabla f(x,y) = \begin{pmatrix} 2x \\ 4/3y^{1/3} \end{pmatrix}$

and

$$\nabla^2 f(x,y) = \begin{bmatrix} 2 & 0 \\ 0 & \frac{4}{9} y^{-2/3} \end{bmatrix}$$

Clearly the minimum is at (x,y) = (0,0) and in this case

$$\nabla^{2} f(x,y)^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & \frac{4}{9} y^{2/3} \end{bmatrix} \xrightarrow{y \to 0} \begin{bmatrix} 1/2 & 0 \\ 0 & 0 \end{bmatrix}.$$

In this case, any algorithm which attempts to generate $\nabla^2 f(0,0)^{-1}$ can be expected to become singular.

In the usual formulation of the Fletcher-Powell algorithm the matrix H_{Γ} is tested at each iteration and if any elements become too large or too small the algorithm is re-started. Thus, in the worse case, the Fletcher-Powell algorithm would de-generate into the usual gradient algorithm. Davidon's algorithm, both the original and as modified here, provide assurance (through α,β) that the computations can proceed by making conservative estimates to the inverse of the Hessian.

The relative speed of Davidon's algorithm compared to conjugate gradient and Fletcher-Powell algorithms was reported on in Reference 11. Those results indicate the Davidon method is the superior algorithm for the solution of most of the test problems considered. In the few cases where another algorithm was found superior the differences were either marginal or, as in a few cases, the Fletcher-Powell algorithm arrived at an exact solution.

SECTION IV

DIRECT APPLICATION OF RANK-ONE

Although rank-one methods have been used to construct algorithms for minimizing arbitrary functions, the direct application of rank-one to minimizing special classes of functions has been overlooked. Consider, for example, minimizing the function

$$f(x) = 1/2||x||^2 + 1/2k(a - m^*x)^2$$

where x, $m \in \mathbb{R}^n$ and k, a are scalars. The gradient of f at x is

$$\nabla f(x) = x + k(a - m'x)(-m).$$

The necessary condition for f to be minimum at $x = x^*$ is that the gradient ($\nabla f(x^*)$) be zero. Furthermore, since m'x is a scalar, (m'x)m can be rewritten as $m(m^*x) = (mm^*)x$, where mm' is the outer or tensor product, i.e., an nxn rank-one matrix. Use this fact and set the gradient to zero to obtain

$$x^* - kam + kmm'x^* = 0$$
or
 $(I + kmm')x^* = kam$

The rank-one matrix inversion technique (Section III) is applied to obtain directly

$$x* = (I + kmm')^{-1}kam$$

$$= (I - \frac{kmm'}{1 + km'm})kam$$

$$= \frac{ka}{1 + km'm}m.$$

More complex forms can be handled with little difficulty. Let M be a m \times n matrix, a an m-vector and consider minimizing

$$f(x) = 1/2 ||x||_{R^n}^2 + 1/2k||a - Mx||_{R^m}^2$$
$$= 1/2x'x + 1/2k(a - Mx)'(a - Mx).$$

The gradient of this function is given by

$$Vf(x) = x + k(-H')(a - Hx) = x - kH'a + kH'Hx.$$

Again, for a minimum of f at x*

$$\nabla f(x^*) = (I + kM^*M)x^* - kM^*a = 0$$

or $x^* = (I + kH'H)^{-1}kH'a$

To clarify the structure of M'M, let m_i be the $i\frac{th}{t}$ row of M represented as a column vector, m_{ij} the $j\frac{th}{t}$ element of m_i as well as the $i,j\frac{th}{t}$ element of M.

Since
$$\{M'M\}_{ij} = \sum_{k=1}^{m} {m_{ki}^{m}}_{kj}$$
,

and

$$\begin{bmatrix}
 m_k m_k^* & = \begin{bmatrix} m_{k1} \\ m_{k2} \\ \dots \\ m_{kn} \end{bmatrix}
 \begin{bmatrix}
 m_{k1}, m_{k2}, \dots, m_{kn} \end{bmatrix}$$

$$\begin{bmatrix} {}^{m}k1^{m}k1 & {}^{m}k1^{m}k2 & \cdots & {}^{m}k1^{m}kn \\ {}^{m}k2^{m}k1 & {}^{m}k2^{m}k1 & \cdots & {}^{m}k2^{m}kn \\ \\ \vdots & \vdots & \ddots & \\ {}^{m}kn^{m}k1 & {}^{m}kn^{m}k2 & \cdots & {}^{m}kn^{m}kn \end{bmatrix}$$

it follows that

$$M'H = \sum_{k=1}^{m} u_k m_k'.$$

Now I + kM*M may be represented as I + $\sum_{i=1}^{m} km_{i}m_{i}^{*}$ and its inverse may be computed recursively by repeated application of the rank-one inversion

method. The question of the existence of $(I + kM'H)^{-1}$ is answered in Theorem II which follows shortly. Formally, the method of computing $(I + kM'H)^{-1}$ is derived recursively as follows:

Set
$$C_0 = I$$
 $D_0 = C_0^{-1} = I$ $C_1 = I + km_1^*m_1^*$ $D_1 = C_1^{-1} = D_0 - \frac{kD_0^*m_1^*m_1^*D_0}{1 + km_1^*D_0^*m_1}$...

 $C_1 = C_{1-1} + km_1^*m_1^*$ $D_1 = C_1^{-1} = D_{1-1} - \frac{kD_{1-1}^*m_1^*m_1^*D_{1-1}^*n_1}{1 + km_1^*D_{1-1}^*m_1}$...

 $C_m = I + kM^*M$ $D_m = C_m^{-1} = (I + kM^*M)^{-1}$

Note D_m is computed in precisely m steps (m = row order of M) and the C_i 's need not be computed since D_i can be computed from D_{i-1} , m_i and k. For this problem, $x^* = D_m(kM^ia)$ is the required solution satisfying $\nabla f(x^*) = 0$.

1. APPLICATION TO AN AIRCRAFT WING-ROOT BENDING PROBLEM

The above rank-one method was used to determine the optimum air-craft wing control surface deflections required to minimize a specified penalty function (References 8,9). The objective of the problem was to reduce wing-root bending moments through active control of trailing edge control surfaces on the wing. In a physical application such reduced bending moment loads could lead to reduced structural requirements, thus reducing the aircraft weight and improving aircraft performance, or alternatively lead to an expansion of the aircraft operational envelope. An immediate consequence of the control surface deflections is to change other important aircraft wing characteristics, principally lift, pitching moment

AFFDE-18-72-77

and drag. These also affect aircraft performance.

For most applications, increased lift at a given angle-of-attack is desirable occause an increase in vertical acceleration can be obtained with less increase in wing angle-of-attack. This leads to less drag and hence more efficient operation.

Any change in wing-generated pitching moment would require changes in aircraft trim to generate balancing moments. Although trim changes can be accomplished automatically, this adds system complexity and interface problems. The desire to avoid complexity is motivated by the usually valid supposition that the more complex system is less reliable.

From drag considerations and because linear aerodynamic theory was used, it was desired to restrict the magnitude of control surface deflections. The use of linear aerodynamic theory, in addition to being conceptually and computationally simpler, was necessary to provide the linear system description.

The mathematical model of this problem was formulated as follows: Changes in wing lift, pitching moment and root bending can be represented, assuming linear aerodynamics, as a linear function of the control surface deflections; therefore let $\delta = (\delta_1, \, \delta_2, \, \ldots, \, \delta_n)$ represent the deflections of the wing control surfaces, $\Delta C_L = m_1^* \delta$ the change in wing lift, $\Delta C_M = m_2^* \delta$ be the change in wing pitching moment, $\Delta C_{RB} = m_3^* \delta$ the change in wing-root bending moment. The following cost function was formulated which when minimized would tend to minimize wing-root bending moment while holding changes in pitching moment small, maximize the change in lift and maintain reasonable control surface deflections:

AFFOL IS-72-TT

$$J(\delta) = K_1(\delta C_{23}) + K_2(\delta C_{11})^2 - K_3(\delta C_{11}) + 1/2[|\delta|]^2$$

where the weighting factors, K_1 , K_2 , and K_3 , are assumed known a priori. The solution to this problem is readily obtained by a single application of the rank-one method.

Although the solutions are valid the results are difficult to interpret due to the fact that each solution represents a different wing-lift condition. The problem was reformulated at a constant lift condition. Specifically, an increment in angle-of-attack was incorporated to allow the wing to generate the same lift with or without controls deflected. The more complete representations for the changes in wing characteristics are:

$$\Delta C_{\underline{L}} = \underline{\mathbf{n}}_{\underline{I}}^{T} \delta + C_{\underline{L}\alpha}^{\alpha},$$

$$\Delta C_{\mathbf{H}} = \mathbf{m}_{2}^{\dagger} \delta + C_{\mathbf{H}\alpha} \alpha,$$

$$\Delta C_{RB} = m_3^* \delta + C_{RB_{cc}} \alpha.$$

Thus for constant lift $\Delta C_L = 0$ and $\alpha = -\frac{1}{C_L} = 1^* \delta$.

Then the changes in pitching moment and wing-root bending at constant lift are:

$$\Delta C_{M} = m_{2}^{2} \delta - \frac{c_{M_{\alpha}}}{c_{L_{\alpha}}} m_{1}^{1} \delta = \overline{m}_{2}^{1} \delta$$

$$\Delta C_{RB} = m_3^{\dagger} \delta - \frac{C_{RB\alpha}}{C_{L\alpha}} m_1^{\dagger} \delta = \overline{m}_3^{\dagger} \delta .$$

Finally, from linear aerodynamic theory, the induced drag of a wing is minimum when the wing pressure distribution is elliptical. Therefore, an additional term was included to represent the deviation of

the predicted pressure distribution from elliptical. Specifically, the spanwise distribution as a function of the normalized semispan η was represented as $c_i(\eta) = c_i(\eta) + c_i(\eta)\alpha + c_i(\eta)\alpha$ and the elliptical distribution as $c_i(\eta) = \frac{2c_i}{\pi} \sqrt{1-\eta^2}$. Following the above procedure for correcting the angle-of-attack change,

$$c_{\mathbf{\hat{\ell}}_{\mathbf{ellip}}}(\eta) - C_{\mathbf{\hat{\ell}}}(\eta) = \left[\frac{2C_{\mathbf{L}}}{\pi}\sqrt{1 - \eta^2} - c_{\mathbf{\hat{\ell}}_{\mathbf{0}}}(\eta)\right] - \left[c_{\mathbf{\hat{\ell}}_{\mathbf{\hat{0}}}}(\hat{\mathbf{e}}) - \frac{c_{\mathbf{\hat{\ell}}_{\mathbf{0}}}(\eta)}{C_{\mathbf{L}_{\mathbf{0}}}}\mathbf{m}_{\mathbf{1}}^{\mathsf{T}}\right]\hat{\mathbf{e}}$$

This is discretized for k spanwise locations $(n_i, i=1, ..., k)$ then a measure of the spanwise distribution deviation from elliptical is given by $||\lambda_0 - 2\delta||^2$, where

$$\{\lambda_0\}_i = \frac{2c}{t} \sqrt{1 - \eta_i^2} - c_t (\eta_i)$$

$$\left\{A\right\}_{i,j} = c_{\ell_{\hat{c}_{i}}}(\eta_{i}) - \frac{c_{\ell}(\eta_{i})}{c_{L_{\alpha}}} m_{lj}$$

and the norm is in Rk.

Finally, the following cost function was formulated to minimize wing-root bending, change in pitching moment, control surface deflections and spanwise distribution error:

$$J(\delta) = K_{1}(\Delta C_{RB}) + \kappa_{2}(\Delta C_{H})^{2} + 1/2||\delta||_{R^{n}}^{2} + K_{5}||\lambda_{o} - \Lambda\delta||_{R^{k}}^{2}$$

$$= K_{1}\overline{m}_{2}^{1}\delta + K_{3}(\overline{m}_{3}^{1}\delta)^{2} + 1/2\delta^{1}\delta + K_{5}(\lambda_{o} - \Lambda\delta)^{1}(\Lambda_{o} - \Lambda\delta).$$

The minimum of this function is obtained by computing the gradient and equating the gradient to zero. The gradient is given by

LFFDL-TB-TZ-TT

$$VJ(\delta) = K_{1}\overline{m}_{2} + K_{3}\overline{m}_{3}\overline{m}_{3}^{*}\delta + \delta + K_{5}(-\Lambda^{*})(\lambda_{0}-\Lambda\delta)$$

$$= (1 + K_{3}\overline{m}_{3}\overline{m}_{3}^{*} + K_{5}\Lambda^{*}\Lambda)\delta + K_{1}\overline{m}_{2} - K_{5}\Lambda^{*}\lambda_{0}$$

and the value of & for which 7J(6) = 0 is given by

$$5 = (1 + K_3 \overline{m}_3 \overline{m}_3' + K_5 A'A)^{-1} (K_5 A'\lambda_0 - K_1 \overline{m}_2).$$

The minimizing & can then be obtained by k+1 iterations of the rankone method.

In both of the above formulations of the problem the weighting factors, K_i 's in the cost function are presumed known. As is often the case in this type problem, this is generally not true. However, the solution to the problem for a given set of K_i 's can be computed extremely rapidly on a digital computer using the rank-one method, making a systematic search over a wide variation of the weighting factors a practical approach to obtaining a realistic solution. For example, a program written for the CDC 6600 computer to solve the above problem computes the basic matrices and vectors required $(h, m_i$'s, etc.), increments the K_i 's in a systematic fashion and computes over 3,000 optimal δ 's in approximately 10.0 seconds central processor, CP, time. With such a volume of data, care had to be taken in the search pattern and output format to avoid an overwhelming output volume.

2. DERIVATION OF THE DIRECT RANK-ONE METHOD

In the preceding development only the solution of the necessary condition for an extremum, $\nabla f(x^*) = 0$, was considered. For the simplest form, a linear combination of quadratic terms with positive coefficients, if an extremum exists it must be a minimum. The addition of linear terms

AFFDL-TB-72-7:

causes no particular difficulty since the quadratic terms can be expected to eventually dominate. For the more general problem where the weighting factors are arbitrary, the rank-one method also provides some information on the mature of the solution.

For any of the above examples the gradient is represented by the general form $\nabla f(x) = Cx + a$, where C is a matrix and a a vector. The second derivative is $\nabla^2 f(x) = C$. The following Lemma provides information on the positive definiteness of C.

Lenna II

If C_n and $D_n = C_n^{-1}$ are positive definite,

$$D_{n+1} = C_{n+1}^{-1} = (C_n + k \mathbf{a}_n \mathbf{a}_n^*)^{-1} = D_n - \frac{k D_n \mathbf{a}_n^* D_n}{1 + k \mathbf{a}_n^* D_n \mathbf{a}_n^*}$$

and 1 + $k \mathbf{m}_n^* \mathbf{D}_n \mathbf{m}_n^* > 0$, then \mathbf{C}_{n+1} and \mathbf{D}_{n+1} are positive definite.

Proof

For an arbitrary $u, u'D_n u > 0$ therefore the ratio

$$\frac{u'D_{n+1}u}{u'D_{n}u} = 1 \cdot \frac{k}{u'D_{n}u} \cdot \frac{u'D_{n}a_{n}^{m}D_{n}u}{1 + ka_{n}^{m}D_{n}a_{n}}$$

$$= 1 \cdot \frac{ka'D_{n}a_{n}}{1 + ka'D_{n}a_{n}} \cdot \frac{(u'D_{n}a_{n})^{2}}{(u'D_{n}u)(n'_{n}D_{n}a_{n})}.$$

By the Cauchy-Schwarz inequality $(u^*D_n m_n)^2 \leq (u^*D_n u) (m_n^*D_m)$ where equality holds if and only if u is a scalar multiple of m_n . Clearly the ratio $\frac{u^*D_n + 1}{u^*D_n u} = 1 \text{ if } u^*D_n m_n = 0 \text{ and this ratio is furthest from one when <math>u$ is a multiple of m_n , in which case

1770L-TR-72-17

$$\frac{\mathbf{u}^{*}D_{m}\mathbf{u}}{\mathbf{u}^{*}D_{m}\mathbf{u}} = 1 - \frac{\mathbf{k}\mathbf{m}_{m}^{*}D_{m}\mathbf{m}}{1 + \mathbf{k}\mathbf{m}_{m}^{*}D_{m}\mathbf{m}}$$

$$= \frac{1}{1 + \mathbf{k}\mathbf{m}_{m}^{*}D_{m}\mathbf{m}}$$

Let $a = \min \{1, (1 + km_n^! D_n m_n)^{-1}\}$

and
$$\beta = \max_{n=0}^{\infty} [1.(1 + km^{*}D_{n}^{m})^{-1}]$$

Then $\alpha u^* D_n u \leq u^* D_{n+1} u \leq \beta u^* D_n u$

Also, for $C_{n+1} = C_n + km_n^n$,

$$\frac{\mathbf{v}^{1}\mathbf{C}_{n+1}\mathbf{v}}{\mathbf{v}^{1}\mathbf{C}_{n}\mathbf{v}} = 1 + \frac{\mathbf{k}(\mathbf{v}^{1}\mathbf{z}_{n})^{2}}{\mathbf{v}^{1}\mathbf{C}_{n}\mathbf{v}}$$

$$= 1 + k v_n^{-1} D_n v_n \frac{(v^* v_n)^2}{(v^* C_n v) (v_n^* D_n v_n)}.$$

Since $C_n = D_n^{-1}$ and C_n and D_n are positive definite

$$(\mathbf{v}^{\dagger}\mathbf{m}_{n})^{2} \leq (\mathbf{v}^{\dagger}\mathbf{C}_{n}\mathbf{v})(\mathbf{m}_{n}^{\dagger}\mathbf{D}_{n}\mathbf{m}_{n}),$$

where equality holds if and only if v is a scalar multiple of m. Thus the ratio $\frac{v^t C_{nv1} v}{v^t C_n v}$ is one if $v^t m_n = 0$ and furthest from one when v is a multiple of m_n in which case

$$\frac{\mathbf{v}^{\dagger} \mathbf{C}_{n+1} \mathbf{v}}{\mathbf{v}^{\dagger} \mathbf{C}_{n} \mathbf{v}} = 1 + k \mathbf{m}_{n}^{\dagger} \mathbf{D}_{n}^{\mathbf{m}}_{n}.$$

Let $\gamma = \min \left\{ 1, 1 + km_{n}^{\dagger} D_{n}^{m} \right\}$

and
$$\lambda = \text{Max} \left\{ 1, 1 + \text{km}_n^* D_n^* m_n \right\}$$

then $\gamma v^* C_n v \leq v^* C_{n+1} v \leq \lambda v^* C_n v$,

Thus the denominator, $1 + km^4D_{nn}$, provides a necessary and sufficient condition for the invertibility of C_{n+1} , i.e., the existence of D_{n+1} , by Lemma I and a sufficient condition for C_{n+1} and D_{n+1} to be positive definite by Lemma II.

Before proceeding to a proof of the rank-one methor, the following observations are made to simplify the notation. For the problem of minimizing

$$f(x) = ||x||_{\mathbb{R}^{n}}^{2} + k_{1}||a - kx||_{\mathbb{R}^{n}}^{2} + k_{2}||\lambda_{0} - \Lambda x||_{\mathbb{R}^{p}}^{2},$$

the gradient is given by

$$\nabla f(x) = 2x + 2k_{1}(-M^{t})(a - Hx)$$

$$+ 2k_{2}(-\Lambda^{t})(\lambda_{0} - Ax)$$

$$= 2(I + k_{1}M^{t}M + k_{2}\Lambda^{t}A)x - 2k_{1}M^{t}a - 2k_{2}\Lambda^{t}\lambda_{0}$$

If x^* minimizes f then $\nabla f(x^*) = 0$ and if $(I + k_1 M^! M + k_2 \Lambda^! \Lambda)^{-1}$ exists, then $x^* = (I + k_1 M^! M + k_2 \Lambda^! \Lambda)^{-1} (k_1 M^! a + k_2 \Lambda^! \lambda_0)$. Furthermore, if $(I + k_1 M^! M + k_2 \Lambda^! \Lambda)$ is positive definite then f has a local minimum at x^* .

As before let m_i , i=1,2,...,m, be the $i^{\frac{th}{t}}$ row of M and λ_i , i=1, ..., p, be the $i^{\frac{th}{t}}$ row of Λ . Then $(I+k_1M'M+k_2\Lambda'\Lambda)$ can be written as

$$I + \sum_{i=1}^{N} \eta_i v_i v_i^*$$

where N = m + p

AFFUL-YR-72-77

$$\eta_{i} = \begin{cases} k_{1} & i = 1, \dots, m \\ k_{2} & i = m+1, \dots, m+p \end{cases}$$

$$\forall_{i} = \begin{cases} n_{i} & i = 1, \dots, m \\ \\ \lambda_{i-m} & i = m+1, \dots, m+p \end{cases}$$

Define C_k and D_k recursively for k = 0,1,...,N as follows:

$$C_{o} = I D_{o} = C_{o}^{-1} = I$$

Theorem II

With $\mathbf{C_k}, \mathbf{d_k}$ and $\mathbf{D_k}$ as defined above

(a)
$$C_{ij} = I + \sum_{i=1}^{N} \eta_{i} v_{i} v_{i}^{i}$$

- (b) so long as $d_{i} \neq 0$, i = 1,...,k $D_{k} = C_{k}^{-1} \text{ and in particular for } k = N,$ $D_{N} = C_{N}^{-1} = (I + \sum_{i=1}^{N} n_{i} v_{i} v_{i}^{*})^{-1}$
- (c) if $d_k > 0$ for all k = 1, 2, ..., N, D_N is positive definite.

Proof

(a) By construction:

$$c_0 = I$$
,
 $c_1 = I + n_1 v_1 v_1'$,

$$C_{2} = C_{1} + \eta_{2}v_{2}v_{2}^{!}$$

$$= I + \eta_{1}v_{1}v_{1}^{!} + \eta_{2}v_{2}v_{2}^{!}$$

$$= I + \sum_{i=1}^{2} \eta_{i}v_{i}v_{i}^{!},$$
...
$$C_{k} = I + \sum_{i=1}^{k} \eta_{i}v_{i}v_{i}^{!}$$
...
$$C_{N} = I + \sum_{i=1}^{N} \eta_{i}v_{i}v_{i}^{!}.$$

(b) By induction:

Initially $C_0 = I$ ar. $D_0 = C_0^{-1} = I$. Suppose the assertion holds for k-1, i.e., $D_{k-1} = C_{k-1}^{-1}$ and $d_i \neq 0$, for $i=1,\ldots,k-1$. Now if $d_k \neq 0$, D_k exists by Lemma I and

$$\begin{aligned} c_k^{D_k} &= (c_{k-1} + \eta_k v_k v_k^{\dagger}) (D_{k-1} - \frac{\eta_k}{d_k} D_{k-1} v_k v_k^{\dagger} D_{k-1}) \\ &= c_{k-1}^{D_{k-1}} - \frac{\eta_k}{d_k} c_{k-1}^{D_{k-1}} v_k v_k^{\dagger} D_{k-1} \\ &+ \eta_k v_k v_k^{\dagger} D_{k-1} - \frac{\eta_k^2}{d_k} v_k v_k^{\dagger} D_{k-1} v_k v_k^{\dagger} D_{k-1}. \end{aligned}$$

Since $C_{k-1}D_{k-1} = I$ and $v_k^{\dagger}D_{k-1}v_k$ is a scalar, the above may be rewritten

as
$$c_k D_k = I - \frac{\eta_k}{d_k} v_k v_k^{\dagger} D_{k-1} + \eta_k v_k v_k^{\dagger} D_{k-1}$$

$$- \frac{\eta_k^2}{d_k} (v_k^{\dagger} D_{k-1} v_k) v_k v_k^{\dagger} D_{k-1}$$

=
$$I - \frac{\eta_k}{d_k} v_k v_k^* D_{k-1} (1 - d_k + \eta_k v_k^* D_{k-1} v_k)$$

= I_*

Where $d_k = 1 + \eta_k v_k' D_{k-1} v_k$ is used. Thus the assertion holds for k. Furthermore, if $d_k \neq 0$ for k=1,...,N then $D_k = C_k^{-1}$ for k=1,...,N and in particular, from (a),

$$D_{N} = C_{N}^{-1} = (I + \sum_{i=1}^{N} \eta_{i} v_{i} v_{i}^{i})^{-1}$$

(c) By induction:

Clearly $D_0 = C_0^{-1} = I$ is positive definite.

Suppose the assertion holds for k-1, i.e., D_{k-1} is positive definite. If $d_k > 0$ then by Lemma II the ratio $\frac{u^*D_{ku}}{u^*D_{k-1}u}$ lies between 1 and

 $d_k^{-1} > 0$ for arbitrary u. Hence, u'D_ku is positive for all u which impositive definite and the assertion holds for k. If $d_k > 0$ for all k = 1, ..., N then D_N is positive definite.

It should be noted that the matrices C_k need not be computed since D_k and d_k can be computed from D_{k-1} , v_k and n_k directly. The denominator d_k provides a convenient check: if d_k is zero (or sufficiently small) then C_k is singular and D_k fails to exist (D_k may be numerically untractable); if d_k is positive at every iteration then the resultant extremum is the minimum since D_N is positive definite. However, if d_k is negative for one or more iterations, D_N may still be positive definite since Lemma II provides only a sufficient condition. In fact it may happen that $I + \sum_{i=1}^{N} n_i v_i v_i$ is invertible, and even positive definite, however an intermediate D_k may fail to exist because d_k is zero. In

such a case, a reordering of the steps will correct the difficulty as shown below.

Suppose all steps up to the $n\frac{th}{n}$, n<N, have been accomplished and $\mathbf{d_n} = 1 - \eta_{\vec{n}} \mathbf{v_n^t} \mathbf{D_{n-1}} \mathbf{v_n} = 0$. Then $\mathbf{D_n}$ fails to exist. The role of $\eta_n \mathbf{v_n} \mathbf{v_n^t}$ and $\eta_{n+1} \mathbf{v_{n+1}} \mathbf{v_{n+1}^t}$ may be interchanged, then $\overline{\mathbf{d_n}} = 1 + \eta_{n+1} \mathbf{v_{n+1}^t} \mathbf{D_{n+1}} \mathbf{v_{n+1}}$ is computed and presumed non-zero. Then

$$\overline{D}_n = D_{n-1} - \frac{\eta_{n+1}}{\overline{d}_n} D_{n-1} v_{n+1} v_{n-1}$$

exists and is the inverse of $\overline{C}_n = C_{n-1} + \eta_{r+1} v_{n+1} v_{n+1}^{\dagger}$

Now the $\text{crm}\ \eta_n v_n v_n^*$ is again considered

$$\overline{d}_{n+1} = 1 + \eta_n v_n^* \overline{D}_n v_n$$

$$= 1 + \eta_n v_n^* (D_{n-1} - \frac{\eta_{n+1}}{d_n} D_{n-1} v_{n+1}^* D_n) v_n$$

$$= 1 + \eta_n v_n^* D_{n-1} v_n - \frac{\eta_n \eta_{n+1}}{d_n} v_n D_{n-1} v_{n+1} v_{n+1}^* D_n v_n$$

$$= - \frac{\eta_n \eta_{n+1}}{\overline{d}_n} (v_n^* D_{n-1} v_{n+1})^2,$$

since the first two terms form $d_n = 0$. In general $d_{n+1} \neq 0$ and hence $D_{n+1} = \overline{D}_n - \frac{\eta_n}{\overline{d}_{n+1}} \overline{D}_n v_n v_n^{\dagger} \overline{D}_n$ exists.

A simple numerical example is presented to illustrate this special case. Consider the problem of inverting $I - \binom{1}{0}(1,0) + \binom{1}{1}(1,1)$. Proceeding as usual, set $D_0 = I$ and compute $d_1 = 1 + (-1)(1,0)(\frac{1}{0}) = 1 - 1 = 0$.

Since D_1 fails to exist, the role of (1,0) and (1,1) are interchanged and \overline{d}_1 is computed as $\overline{d}_1 = 1 + (1)(1,1)(\frac{1}{1}) = 1 + 2 = 3$. Therefore, $\overline{D}_1 = 1 - \frac{1}{3}(\frac{1}{1})(1,1) = \begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix}$

Now, picking up the first term, $\overline{d}_2 = 1 + (-1)(1,0)\overline{D}_1(\frac{1}{0}) = 1 - \frac{2}{3} = \frac{1}{3}$ and $\overline{D}_2 = \overline{D}_1 - \frac{(-1)}{1/3}\overline{D}_1(\frac{1}{0})(1,0)\overline{D}_1$

$$\begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix} + 3 \begin{pmatrix} 2/3 \\ -1/3 \end{pmatrix} (2/3, -1/3)$$

$$\begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix}$$

The correctness of the solution is readily verified as the inverse of

$$1 - {\binom{1}{0}}(1,0) + {\binom{1}{1}}(1,1) = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}.$$

SECTION W

CONCLUSIONS

The main topic of this dissertation, computationa? methods for the solution of unconstrained minimization problems, is covered in three parts; generalized descent algorithms, a rank-one minimization technique for the general problem, and a direct rank-one method for a special class of problems. Although each topic relates to unconstrained minimization, the scope is successively decreased as the specialization is increased.

The theorem on generalized descent algorithms, Theorem I, demonstrates the essential properties of a descent algorithm. Here a descent algorithm is defined to be a computational method in which, at each iteration, a descent direction is generated and a single-dimensional search is conducted for a minimum. Since the setting is highly abstract, additional specialization is included and the theorem is applied to the three descent algorithms in common use: the gradient method, conjugate gradients and the Fletcher-Powell method. The essential property of descent algorithms, choosing a rescent direction and the search for a minimum at each iteration, is sufficient to cause the sequence of function values to monotonically decrease. An additional property is required to assure the derivative of t e function goes to zero.

The descent direction generated by the algorithm must be strict in the sense that convergence of the inner product $[f'(x),\phi(x)]$ to zero must imply convergence of the derivative to zero. This condition is made rigorous in condition (ii) required of the function ϕ . This condition and the uniform continuity of the derivative f' are used to demonstrate the derivative must converge to zero.

Although Theorem I demonstrates the common properties of descent algorithms, it fails to provide any information on one very important property of all such algorithms, convergence rate. Convergence rate involves not only the amount of improvement at each iteration but also, for very practical reasons, the computational time required to accomplish the iterations.

One of the common properties of descent algorithms, the single-dimensional search can also be a drawback in that this is generally the most time consuming step of each iteration. Therefore an algorithm which minimizes a function without requiring repeated single-dimensional searches for a minimum might be superior to any descent method. One such method, Davidon's Variance Algorithm (Reference 13), has been shown to be, in many cases, superior to the three common descent algorithms (Reference 11).

Davidon's method is discussed in Section III.

A new derivation of Davidon's method is presented. This derivation provides a clearer insight into the structure of the algorithm by considering the algorithm as repeated applications of the rank-one matrix inversion technique. Although the rank-one method would be exact when applied to a quadratic form and would require no single-dimensional search, certain added computational precautions must be included if the method is to be used to minimize an arbitrary function. One of these is provided by Davidon and assures the matrix approximation to the inverse of the Hessian is well behaved. Another precaution provided by this author assures only "good" estimates of the solution are used to update the approximation of the inverse of the Hessian. Although a linear search is involved, only the function value is computed and the function need only be decreased not minimized. As a result, for each iteration

both the number of test points is reduced and the complexity of computations at each test point is reduced.

An extension of the Davidon algorithm to the problem of minimizing an arbitrary function defined on an infinite dimensioned space such as a function space is highly interesting. In the finite dimensional case the outer product of vectors, uv', can be represented as a matrix. If u and v belong to a function space, the representation of uv' is not entirely obvious. However, since most problems require digital computation, and the representation of a function is necessarily discretized, the function space may be considered as RN with N very large.

The third area considered in this dissertation applies the rankone method to a special class of problems. Although the class of problems to which the method applies is specialized, it is not uncommon.

In fact many preliminary engineering problems are in this class where
the cost function is a combination of linear and quadratic terms each
weighted by a penalty factor which is constant but unknown. The rankone method provides a rapid solution to the problem. The method also
provides necessary and sufficient tests for the existence of an extremum and a sufficient test for the solution to be a minimum.

Extension of the direct rank-one method to problems in a function space is of continuing interest to the author. As noted for the extension of Davidon's method, care must be observed in the interpretation of the outer product. Since the number of iterations of rank-one required to obtain the solution is determined by the structure of the cost function to be minimized, not the dimension of the underlying vector space, application of the method in a function space has possibilities.

To date the class of problems im a function space which may be solved by the rank-one method is very restricted and the set of known applications is empty. Finally, the problem of representing functions on a digital computer still exists.

To summarize, the following extensions appear to be promising areas for further investigation. Application of Theorem I to other descent algorithms. Application of Davidon's rank-one method to a large sample of test problems to determine in a practical application its computational speed. Application of Davidon's method in a function space might be fruitful if first the direct rank-one method can be successfully applied.

EXCHANGE THE ST

- F. S. Beckman, "The Solution of Linear Equations by the Conjugate Gradient Method", Mathematical Methods for Digital Computers,
 A. Balston and E. S. Wilf (Eds.), John Wiley & Sons, Inc., New York, N. Y., 1960.
- 2. W. C Davidon, Variable Metric Method for Minimization, Argonne National Lab., Report 5990 (Rev.), 1959.
- 3. W. C. Davidon, "Variance Algorithm for Minimization", Computer J., 10 (1968), pp. 406-410.
- 4. R. Fletcher and M. J. D. Powell, "A Rapidly Convergent Descent Method for Minimization", Computer J., 6 (1963), pp. 163-168.
- R. Fletcher and C. M. Beeves, "Function Minimization by Conjugate Gradients", Computer J., 7 (1964), pp. 149-154.
- A. A. Goldstein, "On Steepest Descent", J. SIAM Control, Ser. A, Vol. 3, No. 1, 1965, pp. 147-151.
- 7. A. A. Goldstein, "Minimizing Functionals on Normed Linear Spaces", J. SIAM Control, Vol. 4, No. 1, 1966, pp. 81-89.
- 8. J. E. Jenkins, D. C. Eckholdt, and B. T. Kujawski, "An Assessment of the Interfacing Problem: with CCV Design Concepts", AIAA 2nd Aircraft Design and Operations Meeting, July 1970, AIAA Paper No. 70-926.
- 9. B. T. Kujawski, J. E. Jenkins, and D. C. Eckholdt, "Longitudinal Analysis of Two CCV Design Concepts", AIAA 3rd Aircraft Design and Operations Meeting, July 1971, AIAA Paper No. 71-786.
- 10. L. S. Lasdon, S. M. Mitter, and A. D. Warren, "The Conjugate Gradient Method for Optimal Control Problems", IEEE Trans, on Automatic Control, Vol. AC-12, Nr. 2, April 1967, pp. 132-138.
- 11. Craig E. Miller, A Computational Comparison of Gradient Minimization Algorithms, Masters Thesis. Wright-Patterson Air Force Base, Ohio: Air Force Institute of Technology, March 1971.
- 12. M. J. D. Powell, "A Survey of Numerical Methods for Unconstrained Optimization", SIAM Review, Vol. 12, Nr. 1, (1976) pp. 79-97.
- 13. A. E. Taylor, Advanced Calculus, Ginn and Company, Boston, Mass., 1955.
- 14. M. M. Vainberg, <u>Variational Methods</u> for the Study of Nonlinear Operators, Holden-Day, Inc., 1964.
- 15. H. S. Wilf, "Matrix Inversion by the Annihilation of Rank", J. SIAM, Vol. 7, Nr. 2, June 1959, pp. 149-151.